

Important Dates

- Assignment #1: Intro to Python – Opens 9/13
- Research Proposal Instructions will be handed out. (9/13)
- Assignment #2: Applying Analysis Techniques and Statistical Inference to Data: Opens 10/4
- Assignment #1: Due 10/4
- Students Research Proposals Due: (10/4)

Agenda

- Guest Lecture: Vlad Barash from Graphika
- Lecture 3: Data Munging
- Lab:
 - Web Scraping
 - EDA
- Go over Assignment 1.
- Go Over Project Proposal Instructions.

So you want to be a data scientist...



How will you spend your time?



Enterprise Data Analysis and Visualization: An Interview Study

Enterprise Data Analysis and Visualization: An Interview Study

Sean Kandel, Andreas Paepcke, Joseph M. Hellerstein, and Jeffrey Heer

Abstract—Organizations rely on data analysis to model customer engagement, streamline operations, improve production, inform business decisions, and combat fraud. Though numerous analysis and visualization tools have been built to improve the scale and complexity of these tasks, there has been little research on how analysts actually work. To address this gap, we conducted semi-structured interviews with 35 data analysts from 25 organizations across a variety of industries, including healthcare, retail, marketing and finance. Based on our interview data, we characterize the process of industrial data analysis and document how organizational features of an enterprise impact it. We describe recurring pain points, outstanding challenges, and barriers to adoption for visual analysis tools. Finally, we discuss design implications and opportunities for visual analysis research.

Index Terms—Data, analysis, visualization, enterprise.

1 INTRODUCTION
 Organizations gather increasingly large and complex data sets each year. These organizations rely on data analysis to model customer engagement, streamline operations, improve production, inform sales and business decisions, and combat fraud. Within organizations, an increasing number of analysts are emerging, including “business analyst,” “data analyst” and “data scientist”—perform such tasks. These analysts constitute an important and rapidly growing user population for analysis and visualization tools.

Enterprise data analysis is performed mostly within the context of a larger organization. Analysts often work as a part of an analysis team or business unit. Little research has observed how existing infrastructure, analytical tools, and methods and social connections within an organization impact the analysis process within the enterprise. Understanding how these issues shape analytic workflows can inform the design of future tools.

To better understand day-to-day practices of enterprise analysts, we conducted semi-structured interviews with 35 analysts from sectors including healthcare, retail, finance, and social networking. We asked analysts to walk us through the typical tasks they perform, the tools they use, the challenges they encounter, and the organizational context in which analysts take place.

In this paper, we present the results and analysis of these interviews. We find that our responses are well aligned with previous work that finds that different types of skill in typical workflows. We find that these archetypes vary widely in programming proficiency, reliance on information technology (IT) staff and diversity of tasks, and vary less in their reliance on data science and IT staff. We also find that the diversity of data sources, affect the analysis process. We also describe how collaboration takes place between analysts. We find that analysts often reuse resources such as scripts and intermediate data products. In response, we consider possible impediments to sharing and collaboration.

Next we characterize the analysis process described to us by our respondents. This includes **five high-level tasks**: *discovery, wrangling, profiling, modeling and reporting*. We find that discovery, wrangling, profiling, modeling and reporting. We find that discover-

2 RELATED WORK
 Many researchers have studied analysts and their processes within intelligence agencies [5, 18, 24, 25, 30]. This work characterizes intelligence analysts' process, discusses challenges within the process, and identifies potential improvements. A related body of work exists that much overlap in the high-level analytic process of intelligence and enterprise analysts; these analysts often work on different types of data with different analytic goals. However, these different domains share some tasks. For example, enterprise analysts more often perform analysis on structured data than on documents and emails.

Others have studied the tasks and challenges within the analysis process. For example, Russell et al. [31] propose a set of low-level analysis tasks based on students in an Information Visualization class. Their taxonomy largely includes tasks performed by entry-level analysts. In contrast, the other tasks we have identified, Russell et al. [34] characterize high-level sensemaking activities necessary for analysis. We instead identify specific tasks performed by entry-level analysts and tasks performed by more advanced analysts. The tasks proposed in this domain have limited applicability with other tools, lack support for performing ad hoc transformations of data, and typically do not scale to the necessary volume and diversity of data. We find similar issues across multiple domains.

Fink et al. [9] performed an ethnographic study of cyber security analysts. The tools used by analysts in this domain have limited applicability with other tools, lack support for performing ad hoc transformations of data, and typically do not scale to the necessary volume and diversity of data. We find similar issues across multiple domains.

Several researchers have articulated the importance of capturing provenance to manage analytic workflows [2, 11, 12, 15]. Such sys-

• Sean Kandel is with Stanford University, e-mail: skandek@cs.stanford.edu.
 • Andreas Paepcke is with Stanford University, e-mail: paepcke@cs.stanford.edu.
 • Joseph M. Hellerstein is with UC Berkeley, e-mail: jmh@cs.berkeley.edu.
 • Jeffrey Heer is with Stanford University, e-mail: jheer@cs.stanford.edu.
Manuscript received 31 March 2012; accepted 1 August 2012; posted online 14 October 2012. © 2012 IEEE. Personal use of this material is permitted. For information on obtaining reprints of this article, please send e-mail to: [rsg@computer.org](http://rsg.computer.org).

Interview study of 35 analysts

25 companies
 Healthcare
 Retail, Marketing
 Social networking
 Media
 Finance, Insurance

Various titles
 Data analyst
 Data scientist
 Software engineer
 Consultant
 Chief technical officer

Kandel et al. “Enterprise Data Analysis and Visualization: An Interview Study. IEEE Visual Analytics Science & Technology (VAST), 2012
<http://db.cs.berkeley.edu/papers/vast12-interview.pdf>

“I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I’m lucky if I get to do any ‘analysis’ at all...

... Most of the time once you transform the data ... the insights can be scarily obvious.”

“Once you play with the data you realize you made an assumption that is completely wrong. It’s really useful, it’s not just a waste of time, even though you may be banging your head.”

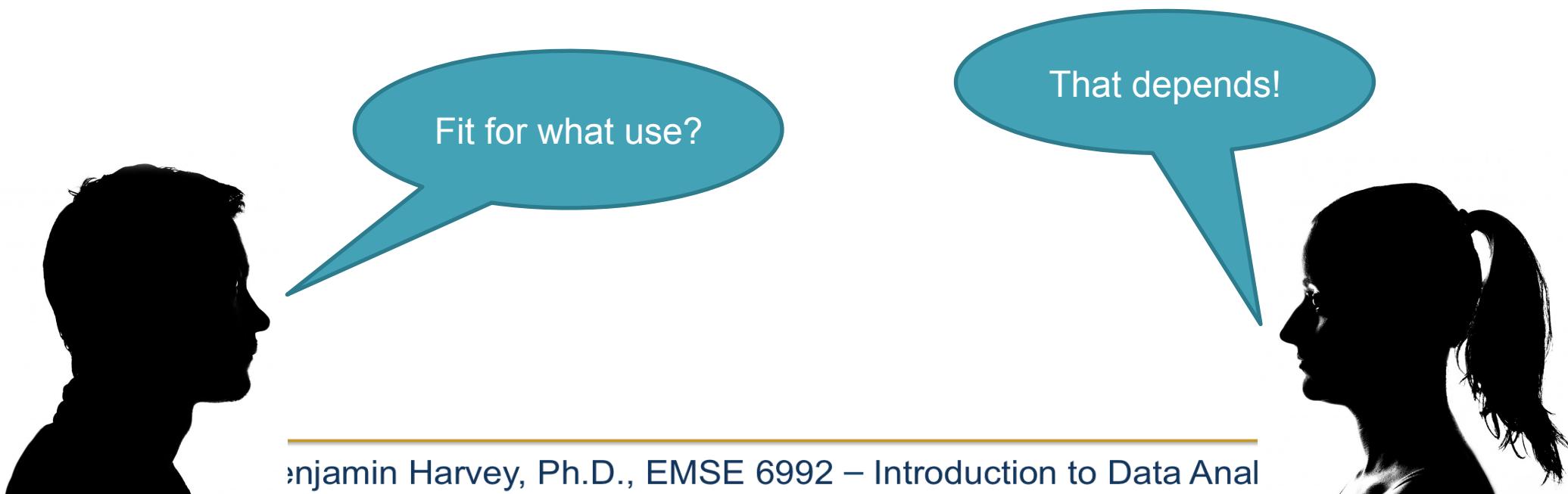
“In practice it tends not to be just data prep, you are learning about the data at the same time, you are learning about what assumptions you can make.”

Data Wrangling



Aka Data Prep, Data Munging, Data Transformation

*Assessing and transforming raw data to make it **fit for use***





Data Wrangling

Aka Data Prep, Data Munging, Data Transformation

Assessing and transforming raw data to make it fit for use

This is how you “get your head in the game”

- Understand what you have
- Assess strengths and weaknesses of your data
- Hypothesize about what to do with your data
- Get it ready

Nobody will know your data as well as you do while wrangling

- Not even the “you” of a few days later

Discussion

When is data “dirty”?

How does that happen?

Today: Data Unboxing and Wrangling

Basic tools:

- UNIX command line
- SublimeText editor

Trifacta: a free visual data wrangling tool

- Born at Berkeley/Stanford
- Codifies some good practices you can also follow “by hand”

Later: Python’s Pandas library

You may choose to use other tools too

- Good data scientists maintain a well-stocked toolbox

A bit of background before we jump in



Data Preparation overview

- ETL
 - We need to **extract** data from the **source(s)**
 - We need to **load** data into the **sink**
 - We need to **transform** data at the source, sink, or in a **staging area**
 - Sources: file, database, event log, web site, HDFS...
 - Sinks: Python, R, SQLite, RDBMS, NoSQL store, files, HDFS...

Data Preparation overview

- Process model
 - The construction of a new data preparation process is done in many phases
 - Data **characterization**
 - Data **cleaning**
 - Data **integration**
 - We must efficiently move data around in space and time
 - Data **transfer**
 - Data **serialization** and **deserialization** (for files or network)

Data Preparation overview

- Workflow
 - The transformation **pipeline** or **workflow** often consists of many steps
 - For example: Unix pipes and filters
 - `$ cat data_science.txt | wc | mail -s "word count" myname@some.com`
 - If the workflow is to be used more than once, it can be **scheduled**
 - Scheduling can be time-based or event-based
 - Use publish-subscribe to register interest (e.g. Twitter feeds)
 - Recording the execution of a workflow is known as capturing **lineage** or **provenance**

The Businessperson

- Data Sources
 - Web pages
 - Excel
- ETL
 - Copy and paste
- Data Warehouse
 - Excel
- Business Intelligence and Analytics
 - Excel functions
 - Excel charts
 - Visual Basic?!

The Programmer

- Data Sources
 - Web scraping, web services API
 - Excel spreadsheet exported as CSV
 - Database queries
- ETL
 - wget, curl, BeautifulSoup, lxml
- Data Warehouse
 - Flat files
- Business Intelligence and Analytics
 - Numpy, Matplotlib, R, Matlab

The Enterprise

- Data Sources
 - Application databases
 - Intranet files
 - Application server log files
- ETL
 - Informatica, IBM DataStage, Ab Initio, Talend
- Data Warehouse
 - Teradata, Oracle, IBM DB2, Microsoft SQL Server
- Business Intelligence and Analytics
 - Business Objects, Cognos, Microstrategy
 - SAS, SPSS, R

The Web Company

- Data Sources
 - Application databases
 - Logs from the services tier
 - Web crawl data
- ETL
 - Flume, Sqoop, Pig, Crunch, Oozie
- Data Warehouse
 - Hadoop/Hive, Spark/Shark
- Business Intelligence and Analytics
 - Custom dashboards: Argus, BirdBrain
 - R

Stages of Wrangling

Raw: Data ingestion & discovery (“unboxing”)

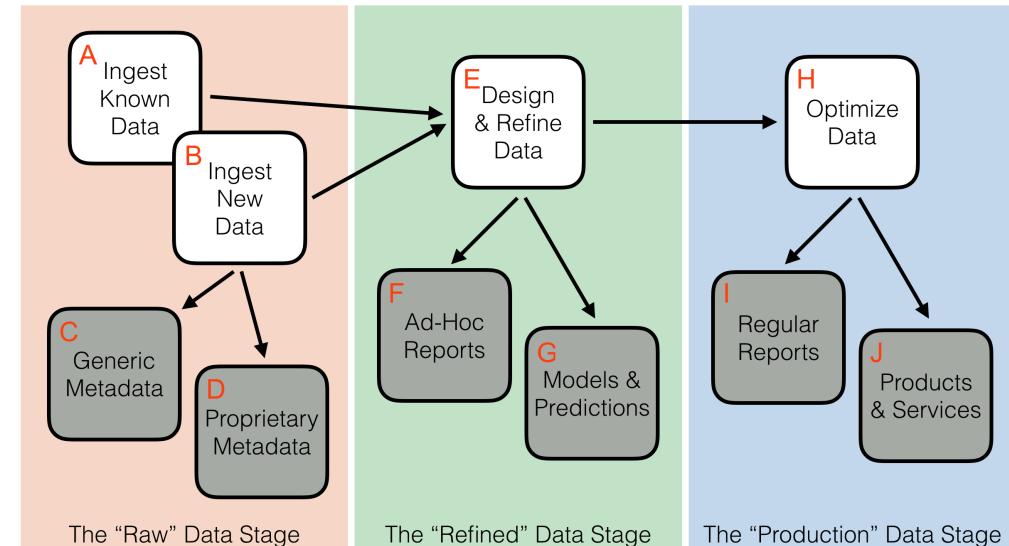
- What: Exploratory *ad hoc* analysis
- Who: The individual wrangler

Refined: Curating data for reuse

- What: Data warehousing, canonical models
- Who: Data curators, IT engineers, actuaries, etc.

Production: Ensuring feeds and workflows

- What: Recurrent, automated use cases:
 - Traditional (e.g. reporting) + New (e.g. recommenders)
- Who: Often involves SW engineers and IT/ops folks



Rattenbury, et al. “Data Wrangling: Techniques and Concepts for Agile Analytics”.
To appear, O’Reilly Media, 2017.

Today

We will focus on the “Raw→Refined” stage

- Unboxing
- Transformation to analytics-ready structure
- Assessment/mitigation of quality issues

This is just a taste

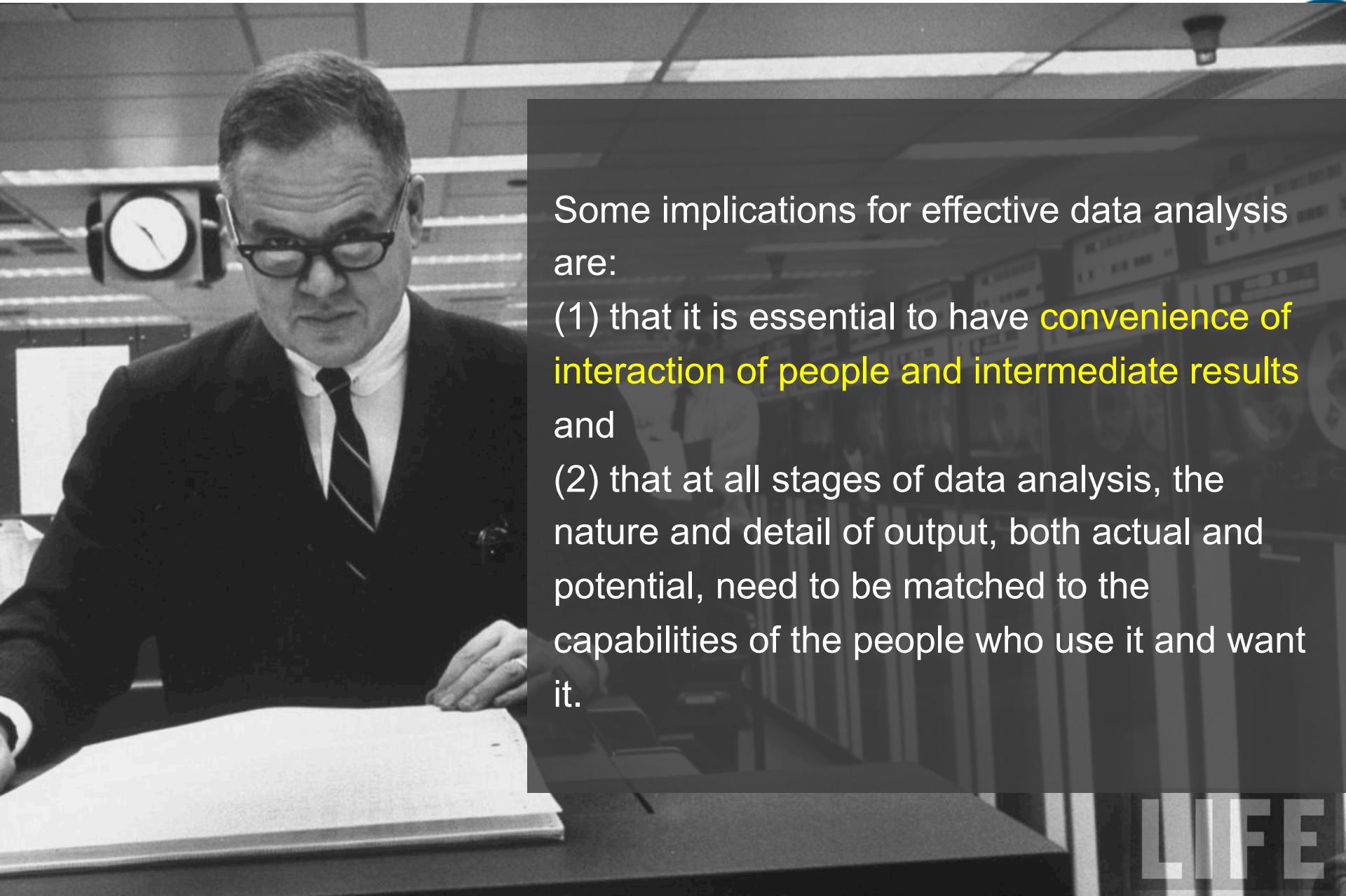
- Soon you will learn to do this in Python & Pandas

More on this in future lectures!



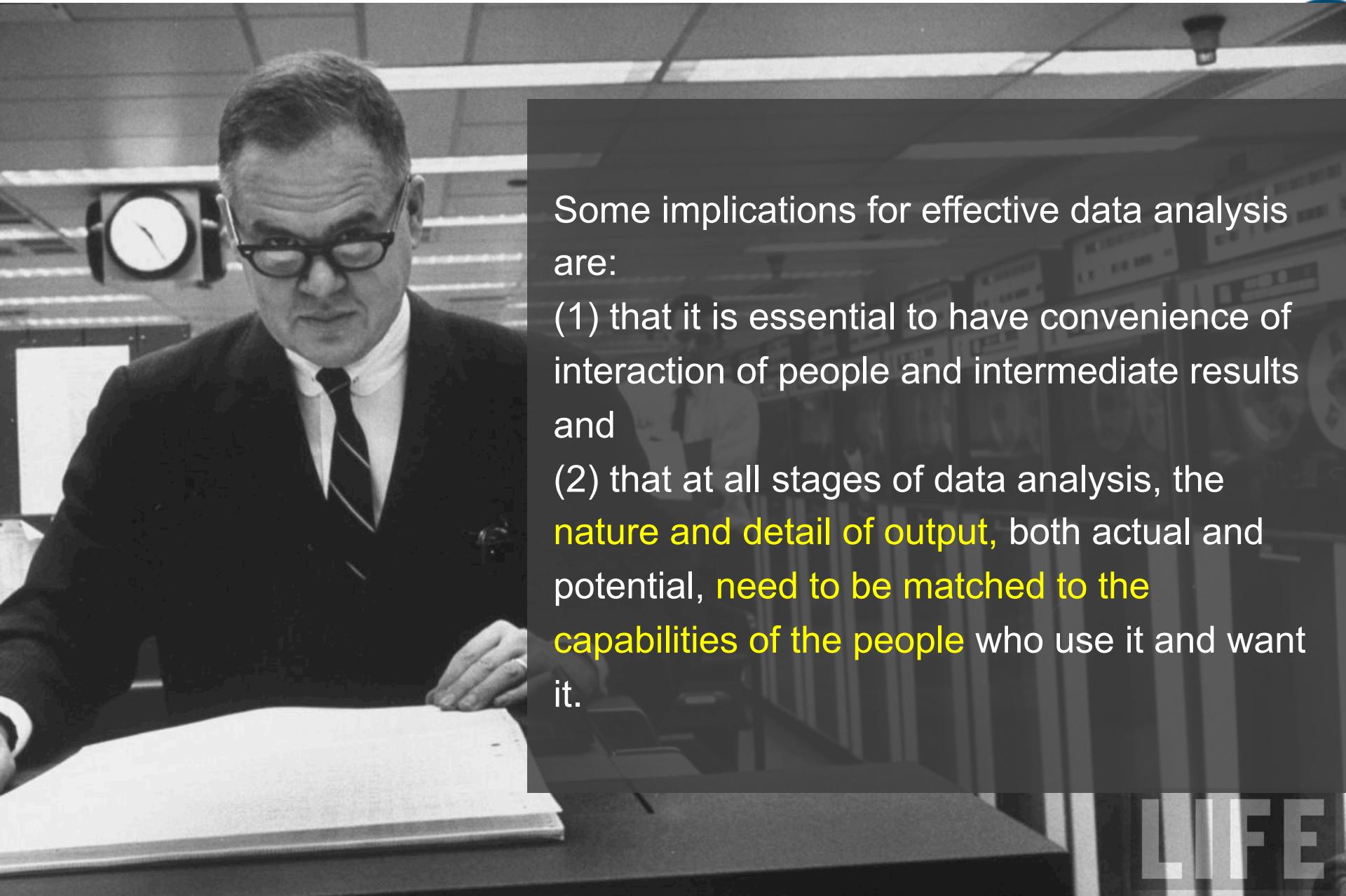
Benjamin Harvey, Ph.D., EMSE 6992 – Introduction to Data Analytics

Data Analysis & Statistics, Tukey 1965



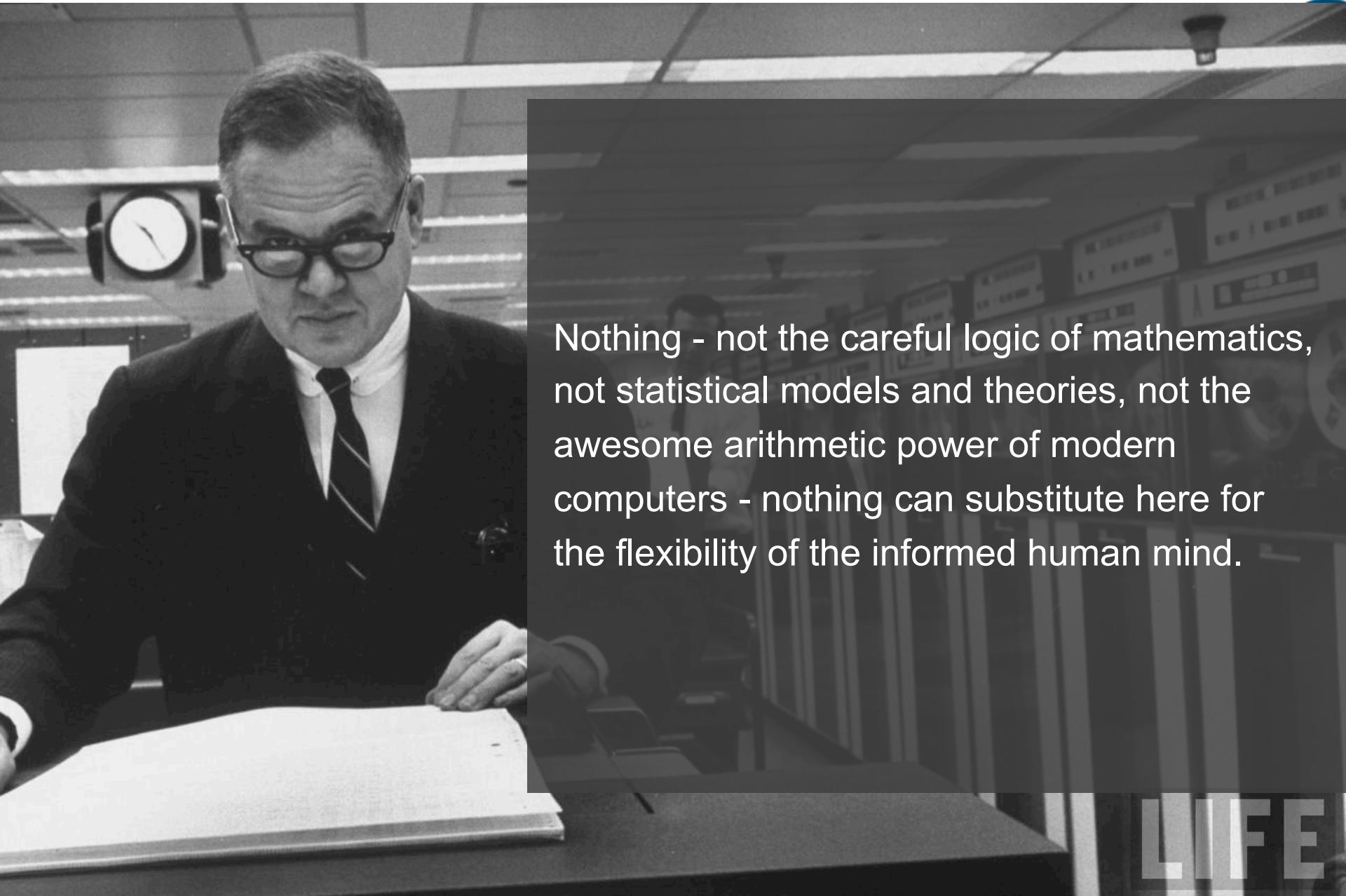
Some implications for effective data analysis are:

- (1) that it is essential to have convenience of interaction of people and intermediate results and
- (2) that at all stages of data analysis, the nature and detail of output, both actual and potential, need to be matched to the capabilities of the people who use it and want it.

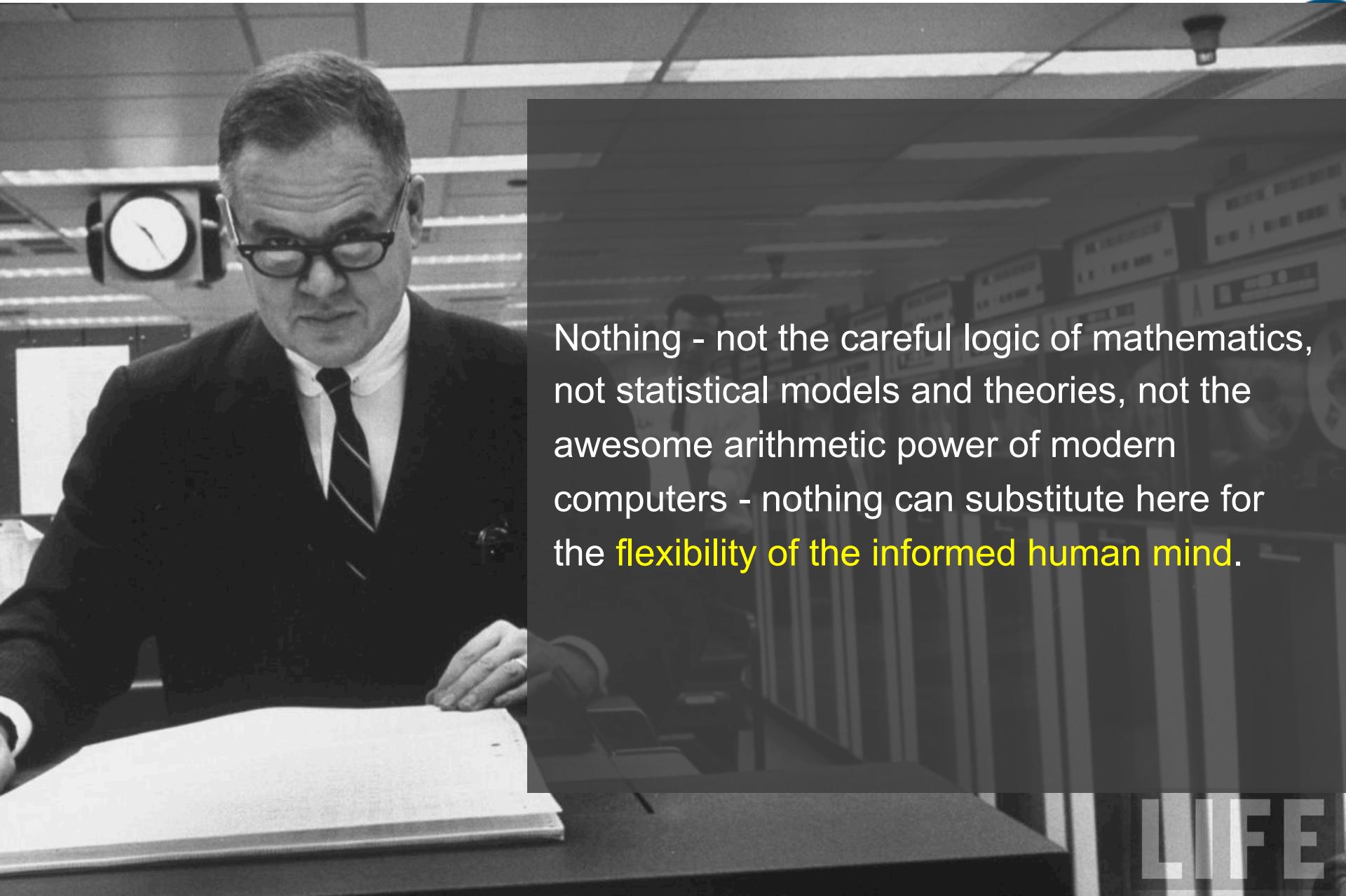


Some implications for effective data analysis are:

- (1) that it is essential to have convenience of interaction of people and intermediate results and
- (2) that at all stages of data analysis, the **nature and detail of output**, both actual and potential, **need to be matched to the capabilities of the people who use it and want it.**



Nothing - not the careful logic of mathematics, not statistical models and theories, not the awesome arithmetic power of modern computers - nothing can substitute here for the flexibility of the informed human mind.



Nothing - not the careful logic of mathematics, not statistical models and theories, not the awesome arithmetic power of modern computers - nothing can substitute here for the **flexibility of the informed human mind**.

Unboxing Data

- What do I have here?
- What do I want to do with it?



These questions rarely have pat answers.

- Typically *contextual* and *user-driven*
- Typically subject to *iterative* cycles of wrangling and analysis

Rough Guide to Wrangling Issues

an outline for today's lecture

- Structure: the “shape” of a data file
- Granularity: how fine/coarse is each datum
- Faithfulness: how well does the data capture “reality”
- Temporality: how is the data situated in time
- Scope: how (in)complete is the data

Many of these are *subjective* qualities! Depend on *context*.

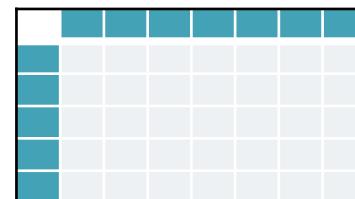
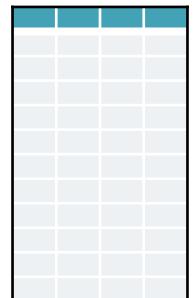
(Data) Science is a human process.

Structure: Rectangular Data

Natural for data analysis: easy to access, filter, tabulate

Two main variants

1. Relations (a.k.a. tables, data-frames)
 - Manipulate with Relational Algebra
2. Matrices
 - Manipulate with Linear Algebra



There are non-rectangular structures as well, e.g. JSON, XML.

For analysis, you will typically need to convert these to rectangular structures.



Unboxing with UNIX Command Line

- File metadata
 - `ls -lh`
 - `file`
 - `WC`
- File (de)compression
 - `gunzip`, `zip`, `bzip`, etc.
- `stdout` and the pipe
- File content:
 - `cat`
 - `head`
 - `tail`
 - `less`
 - `<ctrl>-C`

```
[Lecture4]> file *
D01.gz:      gzip compressed data, was "D01", from Unix, last modified: Mon Jan 23 20:35:45 2017
D02.gz:      gzip compressed data, was "D02", from Unix, last modified: Mon Jan 23 20:36:43 2017
D11.gz:      gzip compressed data, was "D11", from Unix, last modified: Mon Jan 23 20:35:45 2017
D12.gz:      gzip compressed data, was "D12", from Unix, last modified: Mon Jan 23 20:35:45 2017
age classes.txt: ASCII FORTRAN program text
residence area.txt: ASCII text
[Lecture4]>
```

Structure Questions: A Checklist

- Coarse structure
 - Is the data structured as a collection of *records*?
 - How are the individual records delimited in the dataset?
 - How are the record *fields* delimited from one another?
 - Do all records in the dataset contain the same fields?
 - How to access the same fields across records?
By position? Name?
- Are the records nested?
 - No: one atomic (singular) value per field
 - Yes: collection of 0-to-many values in a field
- Encoding
 - How are values encoded? Strings? Codes? Binary?
 - How complex are the individual values?
 - Primitive: numbers and short strings?
 - Unstructured data: natural language text, audio, video
 - Can you decode?

Lecture4] > gunzip -c D01.gz | less

```

<A4><E9><B4><C1>;<B7>;<AD><FB><A5>d<B8><B9>;<A6>~<C4><D6>;<B0>~<EC>;<B0>ö~<A4><C0><C3><FE>;<B0>ö~<BD>s
<BD>X;<BC>zq;<A6><A8><A5><BB>;<BE>P<B0><E2>
2001-01-10 00:00:00;01905168 ;A ;E ;110516;4710063121173;1;20;25
2001-01-15 00:00:00;01637755 ;G ;H ;130101;47106710 ;1;21;26
2001-01-15 00:00:00;01818802 ;C ;E ;320403;20566609 ;1;1062;1395
2001-01-15 00:00:00;01868609 ;G ;C ;100431;4710960748077;1;80;109
2001-01-15 00:00:00;02051994 ;G ;E ;530104;0300410855406;1;112;149
2001-01-15 00:00:00;01786439 ;G ;H ;520470;4711713380322;1;185;269
2001-01-15 00:00:00;00864855 ;G ;D ;100322;4710174041698;1;223;239
2001-01-15 00:00:00;01215854 ;E ;E ;530302;4710084225676;1;108;115
2001-01-15 00:00:00;01779790 ;I ;C ;100403;4901201101014;1;152;189
2001-01-15 00:00:00;01961676 ;E ;F ;500208;4710104104066;1;51;69
2001-01-15 00:00:00;00081269 ;E ;E ;110119;4711258001256;1;30;36
2001-01-15 00:00:00;01403589 ;E ;E ;300301;4710049000386;1;33;45
2001-01-15 00:00:00;00032827 ;F ;E ;130207;4710105011011;1;22;25
2001-01-15 00:00:00;01767322 ;G ;D ;110108;4710088620156;2;154;147
2001-01-15 00:00:00;02125046 ;B ;F ;570206;2210001001236;1;43;54
2001-01-15 00:00:00;01574852 ;C ;F ;321101;20457877 ;1;33;49
2001-01-15 00:00:00;00845298 ;I ;E ;110122;0078895770025;1;54;75
2001-01-15 00:00:00;01658385 ;D ;C ;100324;0723125485025;1;50;65
2001-01-15 00:00:00;02053820 ;D ;F ;760140;4710177041152;1;90;120
2001-01-15 00:00:00;00768023 ;C ;G ;520508;4710017051075;1;52;78
2001-01-15 00:00:00;01669671 ;D ;F ;501102;20008932 ;1;80;112
:

```

Modern Text Editor: SublimeText

- Assessing
 - Coloring
 - Minimap
 - Transformation
 - Do not destroy!
 - Names/Versions?

More on this soon!

transactions.txt

Index	Date	Time	Category	Description	Amount
15369	2001-01-31	00:00:00	018152/b	;B;G;530110;4/10/3530/02/1;1;135;1/9	
15370	2001-01-31	00:00:00	02144580	;C;A;530227;0019200021210;1;145;155	
15371	2001-01-31	00:00:00	00018173	;B;F;130401;4710585930673;1;70;62	
15372	2001-01-31	00:00:00	01565874	;D;E;110402;4710088414021;1;42;55	
15373	2001-01-31	00:00:00	01804492	;C;E;130204;4710105015118;1;77;88	
15374	2001-01-31	00:00:00	00303514	;F;E;120103;4710011401135;1;23;29	
15375	2001-01-31	00:00:00	01798937	;D;E;110117;4711001917018;1;90;103	
15376	2001-01-31	00:00:00	01606515	;F;B;110102;4710311010211;1;19;27	
15377	2001-01-31	00:00:00	00523455	;E;D;760514;4713045022123;1;11;16	
15378	2001-01-31	00:00:00	01740417	;F;E;500705;4715747112300;1;290;399	
15379	2001-01-31	00:00:00	02146263	;C;E;120105;4714381003128;2;78;78	
15380	2001-01-31	00:00:00	018911584	;C;G;560402;8712045008539;1;653;660	
15381	2001-01-31	00:00:00	01470680	;F;F;110507;4710010020061;2;42;50	
15382	2001-01-31	00:00:00	01724974	;C;G;100505;4710018004704;1;35;37	
15383	2001-01-31	00:00:00	01428926	;D;F;100444;4712959001491;1;60;85	
15384	2001-01-31	00:00:00	01283709	;D;F;470103;49099781112950;1;1428;1800	
15385	2001-01-31	00:00:00	01650310	;F;E;711409;4719864068014;2;876;956	
15386	2001-01-31	00:00:00	02072326	;E;H;500202;4710114362029;1;100;107	
15387	2001-01-31	00:00:00	01233621	ré';;10Yd.; ;ÄÖ;USA;°í;°ó~ää;°ó~ä\$éX;¾¶q; ;Ý»,¾Pºå	
15388	2001-02-06	00:00:00	01233621	;B;336;F;110603;4710254011719;2;178;216	
15389	2001-02-11	00:00:00	01750720	;D;502;E;560104;4711553166735;2;86;78	
15390	2001-02-12	00:00:00	00928076	;B;280;E;100201;4712019100607;1;35;49	
15391	2001-02-12	00:00:00	01241525	;E;415;E;130106;4710583110015;1;18;29	
15392	2001-02-12	00:00:00	01667998	;F;679;C;110507;4710740500048;1;120;139	
15393	2001-02-11	00:00:00	01627824	;F;278;F;320506;2160012009601;1;101;155	
15394	2001-02-12	00:00:00	01742800	;H;428;F;100504;4710085104116;2;166;182	
15395	2001-02-12	00:00:00	02125211	;F;251;F;300424;50114080831409;1;102;139	
15396	2001-02-11	00:00:00	01852554	;C;525;F;500705;4715747122262;2;270;208	
15397	2001-02-11	00:00:00	01979510	;E;795;E;110108;4710311703014;1;78;66	
15398	2001-02-11	00:00:00	01517026	;D;170;C;110108;4710395340020;1;70;75	
15399	2001-02-11	00:00:00	01841121	;F;411;D;100201;4712019100607;1;35;49	
15400	2001-02-11	00:00:00	02174983	;J;749;H;530601;4710303267852;1;120;149	
15401	2001-02-12	00:00:00	00197007	;D;970;E;510326;4714232510010;1;58;82	
15402	2001-02-12	00:00:00	02133881	;E;338;G;500202;4710114362029;4;304;328	
15403	2001-02-11	00:00:00	01601473	;D;014;E;500804;9310022862809;3;63;69	
15404	2001-02-11	00:00:00	01852554	;C;525;F;760561;47191111315038;1;10;15	
15405	2001-02-12	00:00:00	01202748	;B;027;F;110507;8801043150606;2;32;42	
15406	2001-02-12	00:00:00	01846683	;K;466;H;740805;4710690099456;1;65;75	
15407	2001-02-11	00:00:00	01080254	;J;802;E;530301;4710094004209;1;50;45	
15408	2001-02-12	00:00:00	01559101	;C;591;E;520464;4711713491783;1;50;70	
15409	2001-02-11	00:00:00	01544206	;H;442;F;530302;4710363502009;2;62;78	
15410	2001-02-11	00:00:00	01663655	;J;636;E;120203;4713758146631;1;90;99	
15411	2001-02-11	00:00:00	01773651	;J;736;A;100431;4712067899218;1;90;109	
15412	2001-02-11	00:00:00	02074788	;C;747;F;100431;4711568404020;1;80;105	
15413	2001-02-12	00:00:00	01648362	;A;483;D;500804;9310022707407;1;21;23	
15414	2001-02-13	00:00:00	01560367	;D;603;F;100205;47100335369510;1;73;89	
15415	2001-02-12	00:00:00	00812931	;D;129;F;530101;4710054132102;2;80;104	

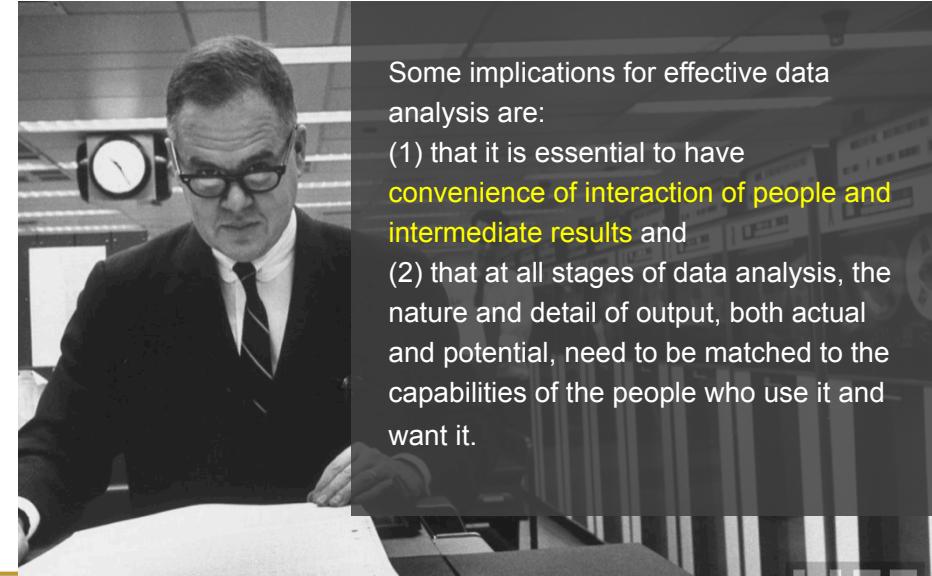
Structure Questions: A Checklist

- Coarse structure
 - Is the data structured as a collection of records?
 - How are the individual records delimited in the dataset?
 - How are the record *fields* delimited from one another?
 - Do all records in the dataset contain the same fields?
 - How to access the same fields across records?
By position? Name?
 - Is the data nested?
 - No: one atomic (singular) value per field
 - Yes: collection of 0-to-many values in a field
 - Encoding
 - How are values encoded? Strings? Codes? Binary?
 - How complex are the individual values?
 - Primitive: numbers and short strings?
 - Unstructured data: natural language text, audio, video
 - Can you decode?

Unboxing with Trifacta

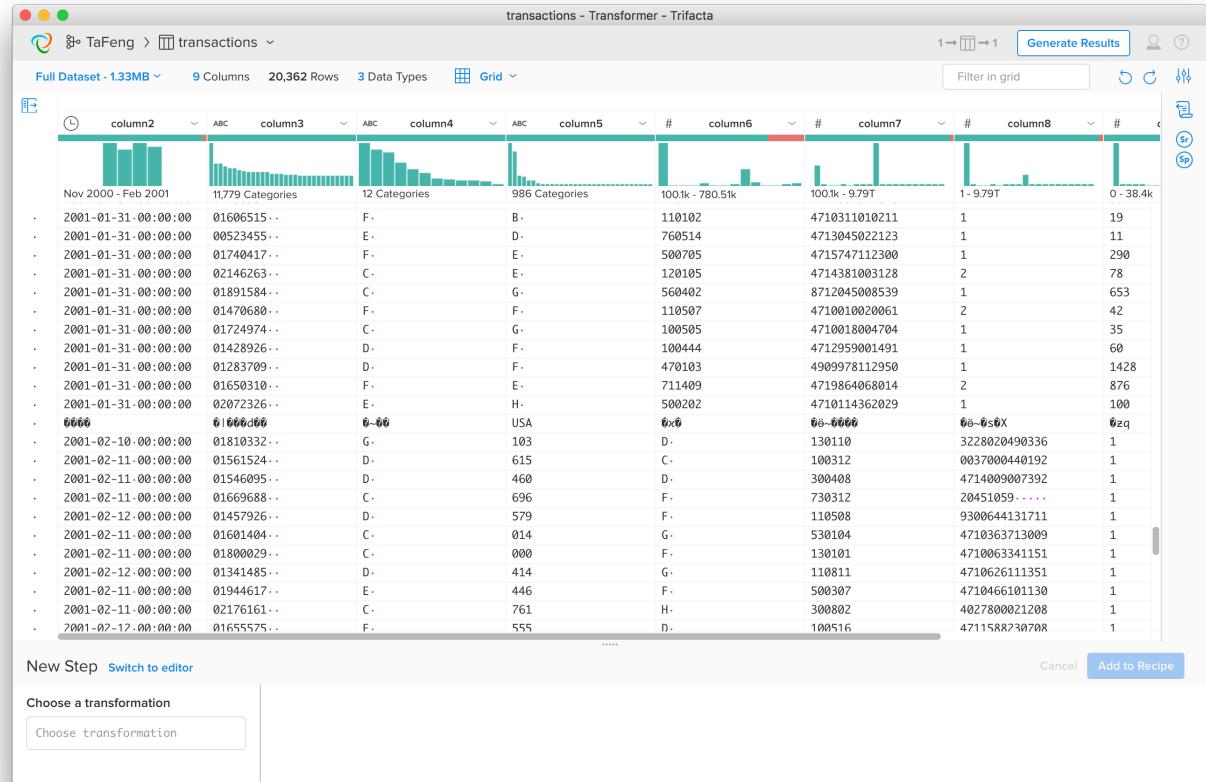


- Visual Profiling + Transformation
 - Iterative and ongoing
- Predictive Interaction
- Transformation language: Wrangle
 - Superset of relational algebra
 - Wrangle (and Pandas) both related to SQL
 - Generates a “recipe” (script)
 - You can run it on data to generate new data
- Scale & Interoperability
 - Paid version works with “big data”
 - Spark/Hadoop



Structure Questions: A Checklist

- Coarse structure
 - Is the data structured as a collection of records?
 - How are the individual records delimited in the dataset?
 - How are the record *fields* delimited from one another?
 - Do all records in the dataset contain the same fields?
 - How to access the same fields across records?
By position? Name?
- Is the data nested?
 - No: one atomic (singular) value per field
 - Yes: collection of 0-to-many values in a field
- Encoding
 - How are values encoded? Strings? Codes? Binary?
 - How complex are the individual values?
 - Primitive: numbers and short strings?
 - Unstructured data: natural language text, audio, video
 - Can you decode?



So Far: Assessing Structure

- Basic tools
 - UNIX commands
 - Text editors
- Simple data transformation
 - Edit-in-place, character-by-character or Find/Replace
 - Save versions of files
 - Scripted transformations
 - Save scripts, inputs, outputs

Value Encodings: Primitive Types

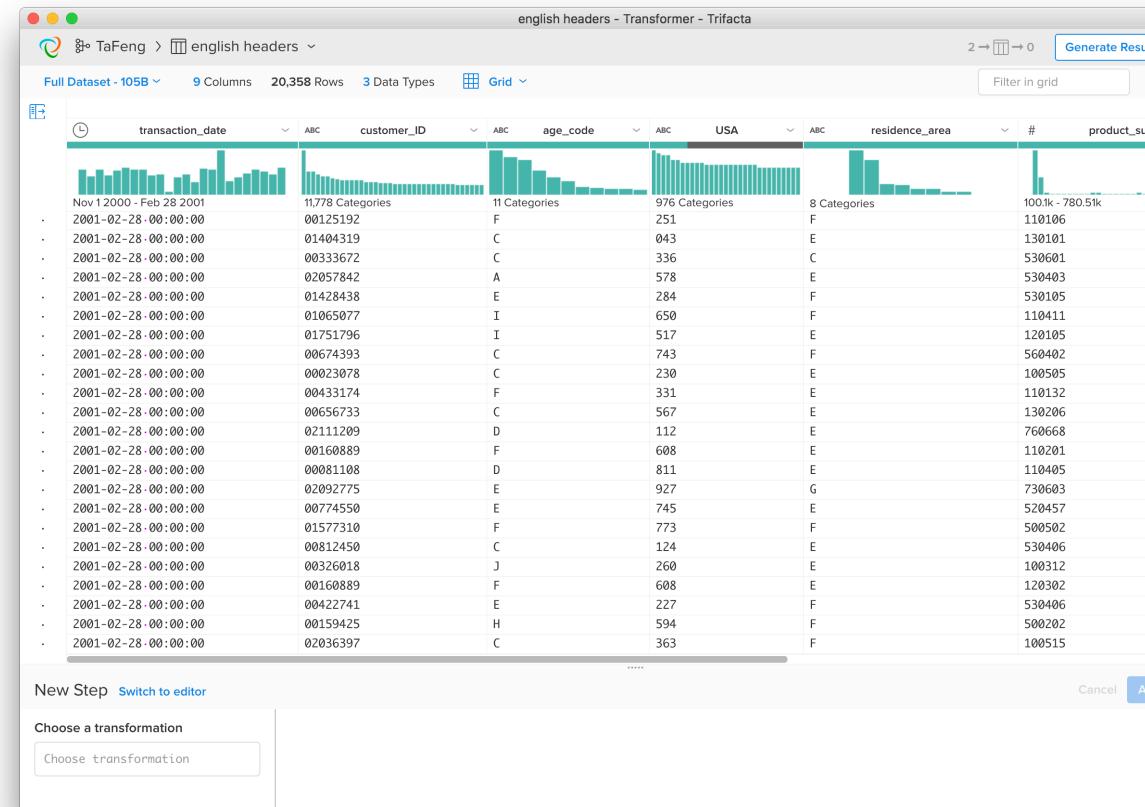
Numbers and Strings

We'll spend time here today

Common special-case “primitives”

- Date/Time
- Geolocations

Space and time can be surprisingly subtle/messy!

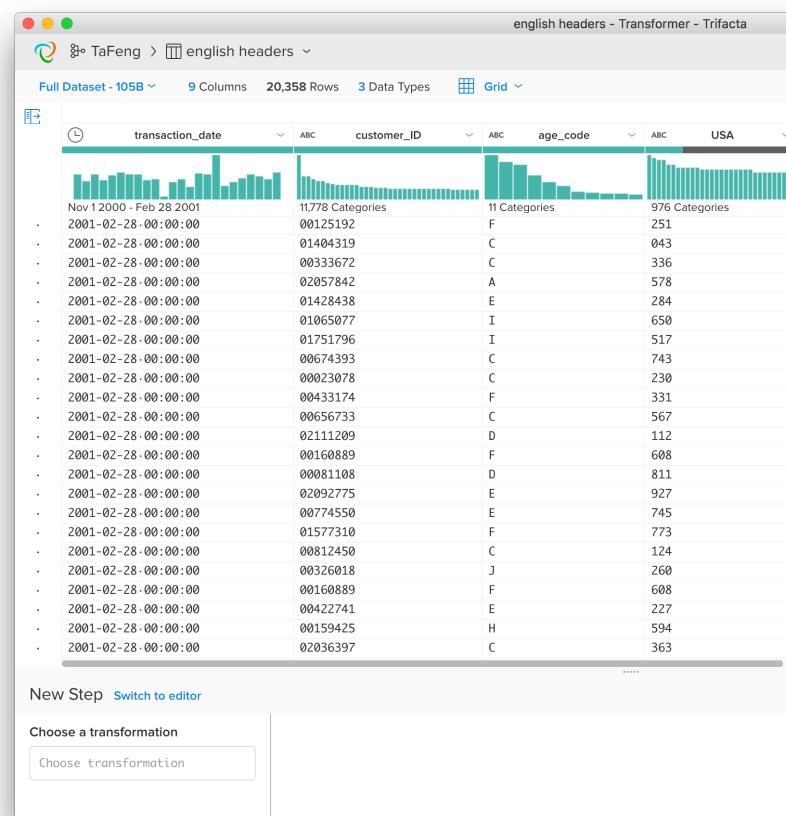


Primitive Type Semantics

- Categorical #1: Nominal
 - E.g. political party affiliation
 - Note: Sometimes even numeric data is nominal
- Categorical #2: Ordinal
 - E.g. education level (none, HS, College, Post-College)
 - Ordered, but not particularly *quantitative*
- Quantitative: amounts, measures
 - Negative or positive
 - Arithmetic “makes sense” (e.g. difference, ratio)
 - Integers vs. Rational numbers

Who cares?

- Affects how you interpret the data
- E.g. In visualization or summarization



In Sum: Tools and Structure

Coarse structure

- Identifying data “shape”, records/fields and delimiters
- Value encodings
- Categorical, Ordinal, Quantitative

Basic tools

- UNIX commands
- Text editors

Simple data transformation

- Edit-in-place, character-by-character or Find/Replace
- Scripted transformations

Assessing Granularity



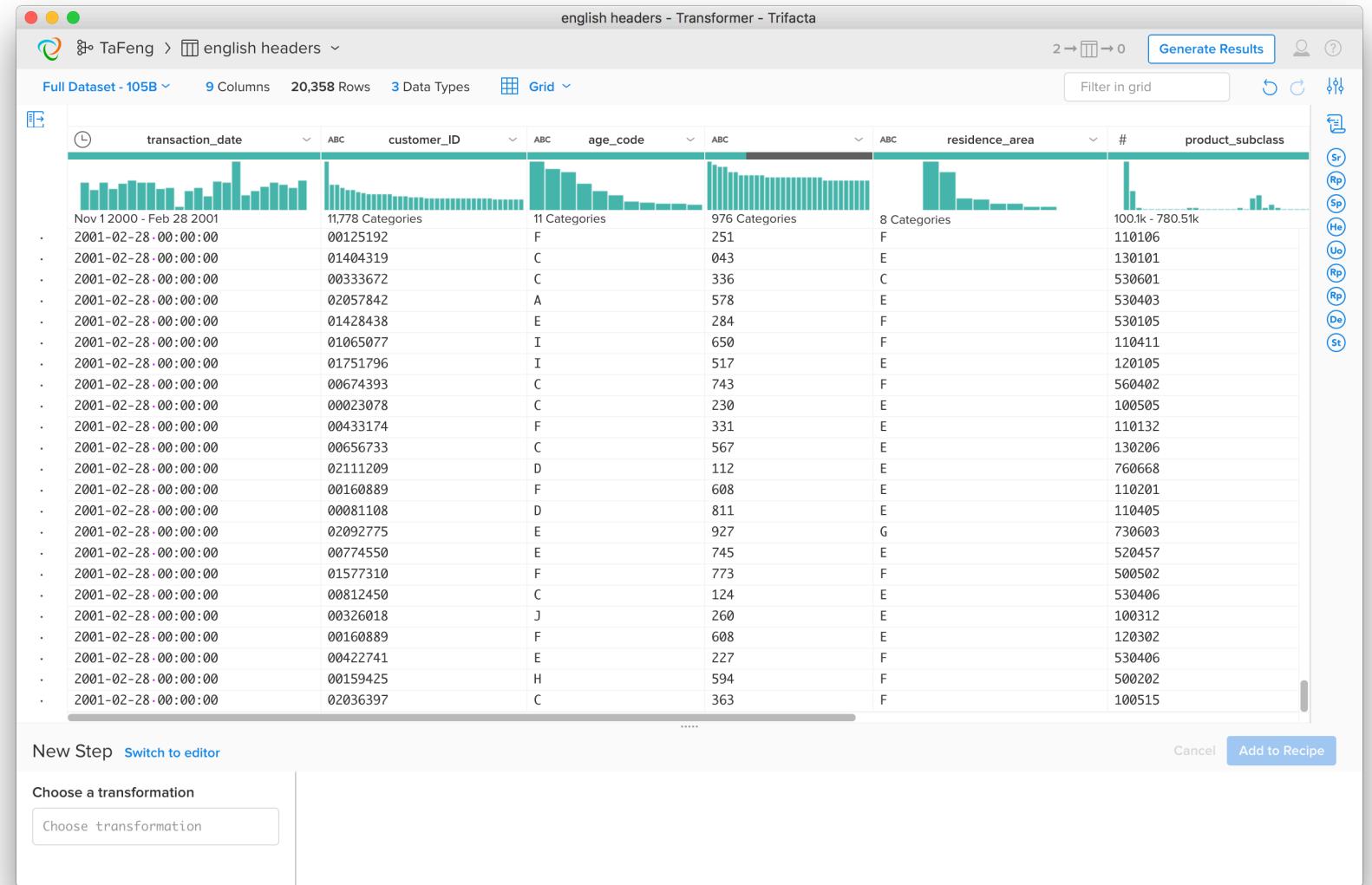
What is the granularity of each record in this data set?

	id	#	favorite_count	source	retweeted
1	763959742755434500	7000			
2	76394972333118700	7942			
3	763949723331185400	5333			
4	763851790329937000	12575			
5	763748672012735600	25779			
6	763558599339191000	20478			
7	763573272954388000	31135			
8	7635196485310000	41672			
9	7635196444100000	18480			
10	763516489655723000	19863			
11	763481802254500000	14104			
12	763474089050000000	19188			
13	763474911831670000	19475			
14	763451452393020000	16823			
15	763451452393020000	20787			
16	76339630123110000	20251			
17	763391459118314000	27981			
18	763385288295054000	17173			
19	763385288295054000	36468			
20	763380978201130000	45699			
21	763167670174252000	18125			
22	763167345405784000	11878			
23	7631195713099981000	23645			
24	763119505629100000	0			
25	763106393056291000	0			
26	763106393020000000	0			
27	7631063833899770000	0			
28	7629823608530505000	17882			
29	7629821427836764000	13658			
30	7629821427836764000	44682			
31	7629714890000000000	11296			
32	7629421269000000000	25510			
33	7628157554235187000	0			
34	7627961673260500000	17211			
35	7627818262649600000	33203			
36	7627818262649600000	25207			
37	762777166476250000	16379			
38	7627764406810952000	15253			
39	762775257565143000	13452			
40	7627743213080679000	11116			
41	7627743213080679000	11983			
42	7627743213080679000	12000			
43	7627439236972200000	0			
44	76269882571988000	21419			
45	7626415954391900000	13946			
46	7625391920112300000	19351			
47	7625371200000000000	26460			
48	7624008608581150000	41271			
49	762399109407065000	0			
50	762848533416750000	33375			
51	7621109182113100000	20572			
52	7621109182113100000	21188			
53	7621044117956110000	8596			
54	7620104261822960000	25958			

TRANSFORM FIDDLER

What does each TaFeng record represent?

Theories?
Justification?



Assessing Granularity



Key (“primary key”)

300 Categories

Given a relation, a key attribute uniquely determines the values in each record

- E.g. an *identifier* like `transaction_id`
- Can be *composite*: e.g. `(City, State)`

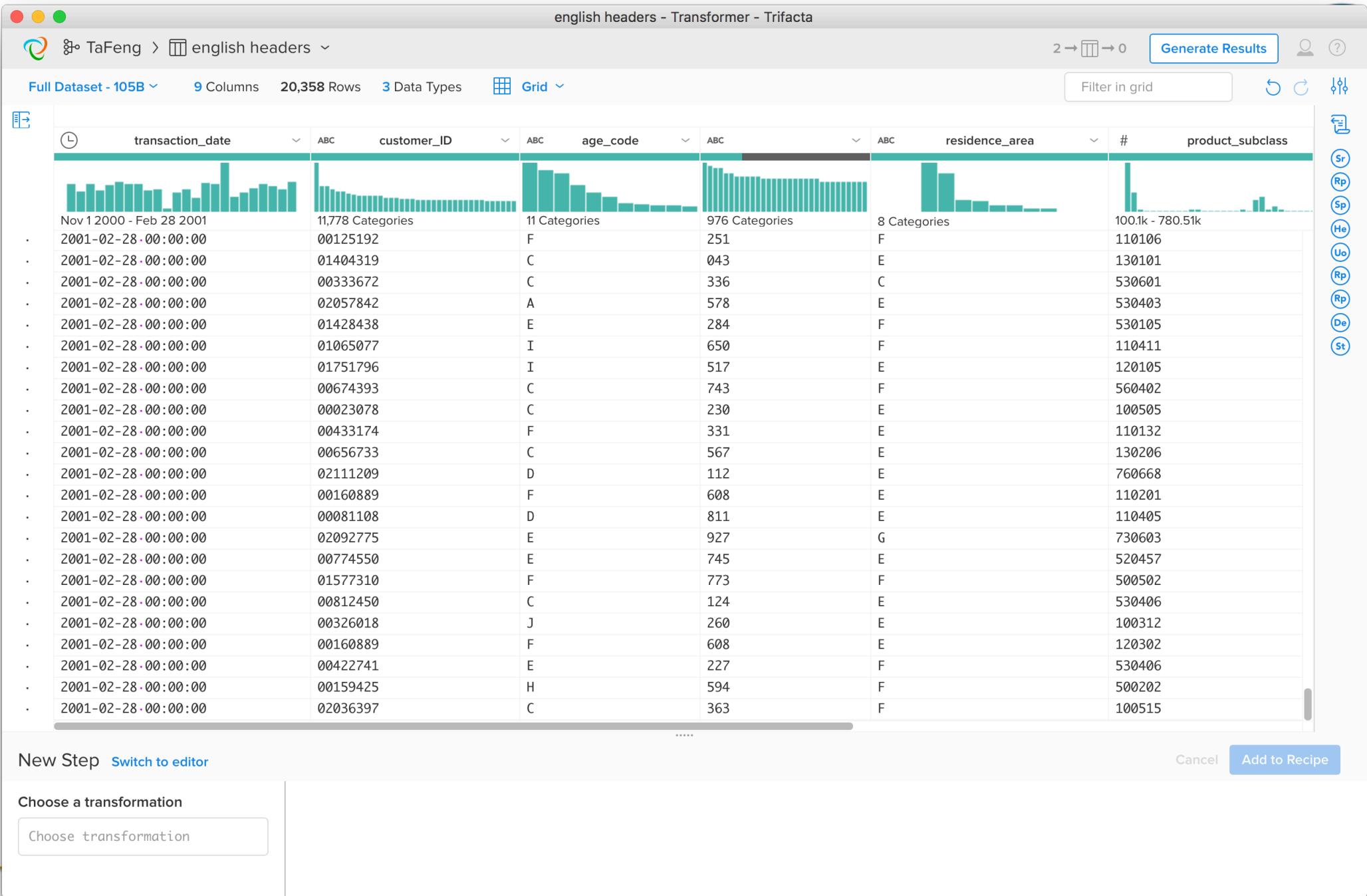
Each value occurs **at most once** in the key column(s).

Hence the semantics (meaning) of the key determines granularity

- `customer_id?` `transaction_time?`
- `(age_code, day_of_week)`

Be careful!

- Goofy key choice? Goofy data.
- Real-world data may have “noisy”, duplicated keys to clean up



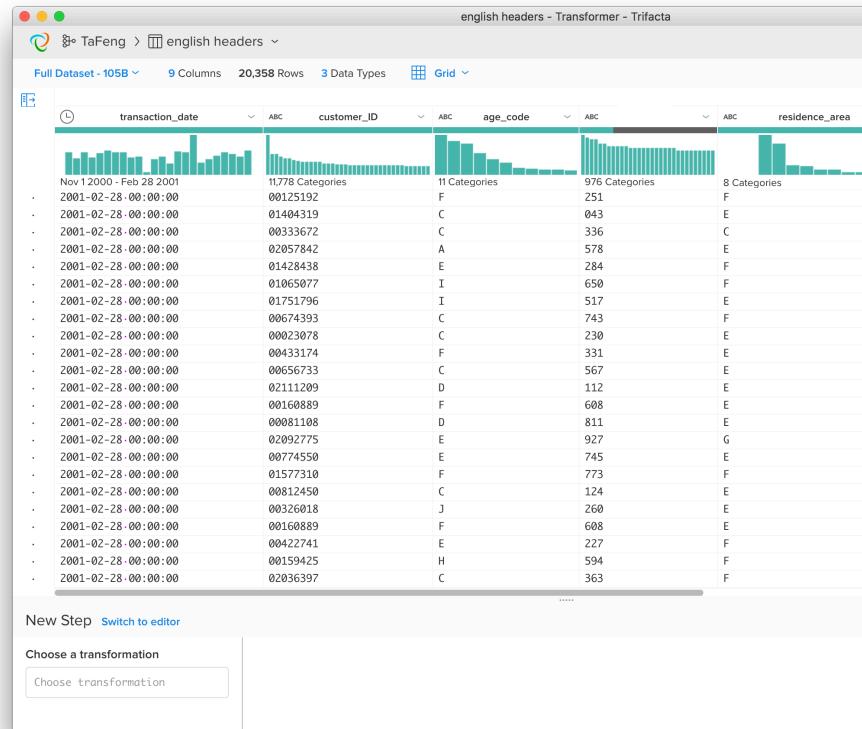
Joining tables with keys

- Keys in our data
 - `age classes.txt`: code
 - `residence area.txt`: code
 - `transactions.txt`: transaction_id
- Transactions.txt also has *foreign keys*
 - i.e. attributes that reference the key of another table
 - `age_code` is a foreign key to `age classes.txt`
 - `residence_area` is a foreign key to `residence area.txt`
- Joining on a foreign key is a “lookup”
 - At most one match for each row of `TaFengTransactions.txt`

More on this in future lectures!

Granularity Questions: a Checklist

- What kind of thing does each record represent?
- Do all records capture granularity at the same level?
- What alternative interpretations of the records are there?
 - E.g. Is this a list of transactions?
Or a list of customers?
- What kinds of aggregation is possible/desirable?
 - From individual people to demographic groups?
 - From individual events to totals across time or regions?
 - Hierarchies (city/county/state, second/minute/hour/day)



So Far: Assessing Granularity

Identifying (primary) keys can help assess Granularity

- Each record is the information *per key*

Foreign Keys are pointers to individual rows in other data sets

- To lookup a row in another table: join the foreign key to its primary key

Assessing Faithfulness

Theme: the Faithfulness of a record can only be evaluated in context

- Application context
- Context in your data set
 - Across records

Students			
id: integer	DOB: date	GPA: float	Risk: float
123457	01/16/1997	3.2	465
123458	01/24/2017	2.7	28
123457	01/16/2002	5.0	27
123459	03/22/1996	3.6	31
123460	06/13/1997	2.2	43

Faithfulness Across Records: Outliers

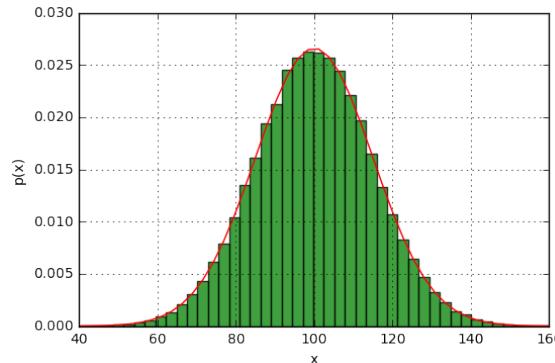
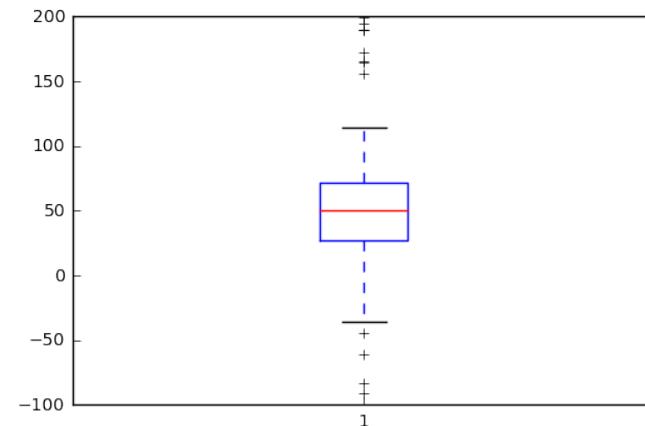
What is an “outlier”?

More on this in future lectures!

- A value that is “far” from the “center”

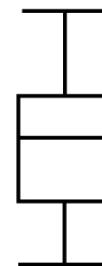
Distribution-based definition

- Center (e.g. average, median)
- Spread (e.g. standard deviation, IQR)



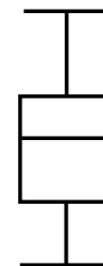
What to do with Outlier

Delete. (“trimming”)



Set to a default

- E.g. the nearest non-outlier (“Winsorizing”)



Good Hygiene:

- Leave the original column
- Derive an indicator column to flag presence of outlier
- Derive a clean column for your use

english headers - Transformer - Trifacta

2 → → 0 [Generate Results](#)

[Full Dataset - 105B](#) 9 Columns 20,358 Rows 3 Data Types [Column Details](#)

Sort: Default [Edit](#)

customer_ID

SUMMARY

	Valid •	20,358	100.0%
Unique	•	11,778	57.9%
Outliers	•	8	0.0%
Mismatched •	•	0	0.0%
Missing •	•	0	0.0%

TOP VALUES

Value	Count
00020459	41
02112589	21
02112596	21
01647457	19
00426053	17
02133874	17
01847994	15
02113579	14
00380393	13
00570565	13
01660562	13
01720006	13
02019604	13
01359015	12

MISMATCHED VALUES

None

OUTLIERS

Value
20002000

STATISTICS

Statistic	Value
Minimum	1,823.00
Lower Quartile	969,543.00
Median	1,595,128.00
Upper Quartile	1,856,156.00
Maximum	20,002,000.00
Average	1,406,004.94
Standard Deviation	712,452.70

VALUE HISTOGRAM

FREQUENT VALUES

New Step [Switch to editor](#)

Choose a transformation

[Choose transformation](#)

[Cancel](#) [Add to Recipe](#)

Assessing Faithfulness Within Records: Dependencies and Correlations



Nov 1 2000 - Feb 28 2001 11,778 Categories 11 Categories 976 Categories 8 Categories 100.1k - 780.51k

2000-11-02 00:00:00	00020459	E		D	130204
2000-11-17 00:00:00	00020459	E		D	100102
2000-11-28 00:00:00	00020459	E		D	110401
2000-11-29 00:00:00	00020459	E		D	100402
2000-11-30 00:00:00	00020459	E		D	500307
2000-12-05 00:00:00	00020459	E		D	100414
2000-12-14 00:00:00	00020459	E		D	760205
2000-12-16 00:00:00	00020459	E		D	120115
2000-12-23 00:00:00	00020459	E		D	100312
2000-12-23 00:00:00	00020459	E		D	120106
2000-12-26 00:00:00	00020459	E		D	100315
2000-12-26 00:00:00	00020459	E		D	100312
2001-01-14 00:00:00	00020459	E		D	730144
2001-01-14 00:00:00	00020459	E		D	110106
2001-01-14 00:00:00	00020459	E		D	100308
2001-01-06 00:00:00	00020459	E		D	500535
2001-01-16 00:00:00	00020459	E		D	110402
2001-01-16 00:00:00	00020459	E		D	110411
2001-01-16 00:00:00	00020459	E		D	110402
2001-01-17 00:00:00	00020459	F		D	100110

SUGGESTIONS

- keep ABC customer_ID 00020459
- delete ABC customer_ID 00020459
- set ABC customer_ID 01172829 ABC customer_ID 01172829
- derive ABC customer_ID false

Cancel Modify Add to Recipe

Functional Dependencies (FDs)

More on this in future!

- Generalization of Keys
- Attribute A *determines* Attribute B
 - $\text{customer_id} \rightarrow \text{age_code}$
 - i.e. $\text{age_code} = f(\text{customer_id})$
 - Function represented as lookup table
- More generally, a set of columns determines another set of columns
 - $\{\text{transaction_time}, \text{customer_id}\} \rightarrow \{\text{age_code}, \text{residence_area}\}$
- Primary Keys are special FDs
 - Right-hand-side is the set of *all* attributes in the relation

Correlations

More on this in future!

Dependence (i.e. lack of independence!) between 2 random variables

Think of the attributes in a relational schema

- An *instance* of that relation was generated from some real-world process
- Each column of that relation is a “random variable” generated by the process

A Functional Dependency is a “deterministic” correlation

Correlations are more general: statistical relationships

- amount and sales_price are correlated

What to do about bad FDs/Correlations

- Cleaning Noisy FDs: $\text{customer_id} \rightarrow \text{age_code}$ (kinda)
 - For a few customer IDs, there are multiple values of age_code
 - Set offending right-hand-side values to all match
 - Set offending left-hand-side values to NULL
- Cleaning Correlations: height correlated with weight
 - Some rows don't seem to follow the correlation (how do we decide?)
 - Can *impute* a likely value for one side or the other
 - Can set one side or the other to NULL
- Don't forget Good Hygiene!
 - Leave the original column
 - Derive new columns (indicators and/or cleaned data)

Careful! More on this in future.

Faithfulness Questions: A Checklist

- Type-specific Faithfulness checks
 - Are dates and addresses legal/reasonable?
 - Numeric codes legal?: E.g. phone numbers, credit cards, social-security numbers etc.
 - Can you validate network endpoints? Email addresses, IP addresses, social network names
 - Need to deduplicate named entities: misspellings, acronyms (UCB vs Berkeley)
- Checks for data entry problems
 - Frequency outliers for common default entry values (00000, 1234567)
 - Misspellings (compare to dictionaries)
 - Sensor drift – often timeseries-based
 - “Curbstoning” in surveys
- Quantitative dirty data
 - Outliers, FDs, Correlations
- Check distribution of inaccuracies
 - Do inaccuracies seem to affect a large fraction of records?
 - Are they concentrated in a particular subset of records?

Tricky!

Summing Up: Faithfulness

- Outliers
- Functional Dependencies & Correlations
- Good Data Cleaning Hygiene
 - Don't overwrite: use indicators and derived columns

Granularity Transformations

- We can coarsen the granularity by picking new keys and “rolling up”
 - GroupBy and Aggregation
- GroupBy
 - Choose a new primary key (the group-by columns)
 - Result will have one row per distinct values of this primary key
- Aggregation
 - Summary (rollup) results per group
 - E.g. count(), or aggregation functions on attributes (sum(x), average(x), stdev(x) etc.)

Finalizing Ta Feng

Full Dataset - 105B ▾ 5 Columns 88 Rows 2 Data Types Grid ▾ Filter in grid

ABC	age_code	ABC	age_range	ABC	residence_area	#	row_count	#	sum_sales_price
	11 Categories		10 Categories		8 Categories				
.	D	35-39	E			5 - 1.57k		447 - 199.04k	
.	D	35-39	F			1573		199038	
.	E	40-44	E			1449		185955	
.	C	30-34	E			1496		171978	
.	C	30-34	F			1226		158129	
.	F	45-49	E			1212		156652	
.	E	40-44	F			1005		133835	
.	D	35-39	G			1098		133341	
.	G	50-54	E			491		120883	
.	F	45-49	F			573		80697	
.	E	40-44	G			685		80573	
.	B	25-29	E			294		78062	
.	B	25-29	F			591		72710	
.	K		H			588		69808	
.	G	50-54	F			412		60434	
.	C	30-34	G			383		57537	
.	E	40-44	C			360		56245	
.						386		56065	
								

New Step [Switch to builder](#)

Enter transform expression [Add to Recipe](#)

What About Non-Rectangular Data

- Nested Data
- Free Text (i.e. for humans)
- Example that has both: Twitter feed

Unnesting Nested Data Types

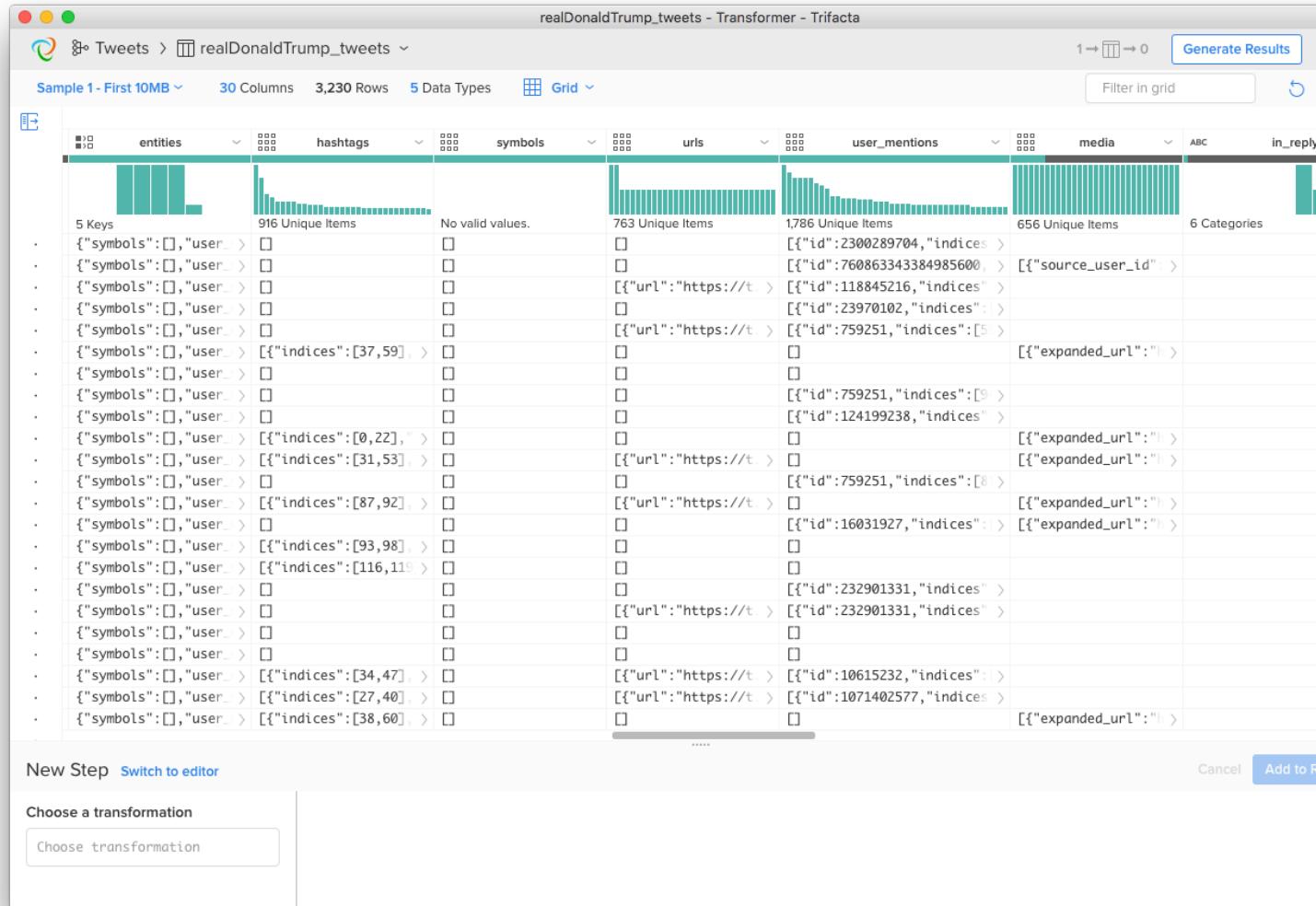
More on this in future lectures!

Maps

- A.k.a. dictionaries, hashes
- A set of **key:val** pairs

Arrays

- Or lists



More on this in future
lectures!

Playing with Text I: String Manipulation

Regular expression (regex) extraction

The screenshot shows the Trifecta Data Transformer application interface. The main window displays a preview of a dataset named "realDonaldTrump_tweets" with 3,230 rows and 27 columns. The "text" column is selected, and a histogram indicates there are 913 unique items. Below the preview, a list of tweets is shown, each containing a sequence of words separated by spaces. At the bottom of the screen, a "New Step" dialog is open, titled "Switch to editor". It contains fields for "Choose a transformation" (set to "extractlist"), "Column" (set to "text"), "On pattern" (set to "\w+"), and "Delimiter" (set to "Edit pattern").

Playing with Text II: Natural Language

Natural Language Processing: a slippery slope

- Entity Resolution
 - IBM vs. International Business machines
- Named Entity Recognition
 - Steve Jobs munched on an apple as he announced new jobs at Apple.
- Sentiment
 - My windows box crashed yet again.
- Etc. Etc.
- Hot topic today: Dialog systems:
 - Alexa, where did I leave my keys?
 - I don't know. It's a pity you don't have a tracker on your key fob.
 - A what?
 - A bluetooth tracker – there are variety of them for sale
 - Stop trying to get me to buy stuff!
 - Sorry about that; I know you've been trying to save this month...
 - ...but if you had one I could locate your keys.

Playing with Text II: Natural Language

- Some simple things
 - Term frequencies
 - Simple Sentiment analysis

Quick Note: Assessing Temporality

- Often two kinds of time in data
 - Time of data entry
 - Time of a recorded phenomenon being “true”
 - E.g. A physical time of an event happening
 - E.g. An “effective” time, e.g. date that a subscription will start
- Often more
- Time is tricky!
 - Periodicities (recurring patterns in Days of the week, or Months of the year)
 - Non-uniform hierarchy of units (# days in a month, # of days in a year, etc.)
 - Time zones are complex: especially daylight savings (summer) time
 - Clocks can be skewed
 - Relativity: true perception of event may vary (yep!)
- Assess timestamps in data carefully!!



Quick Note: Assessing Scope

- Do you have all the data you need
 - Missing columns
 - Join in external data
 - Missing rows/values?
 - Look for sequential patterns with breaks
 - Is this a good sample?
 - Granularity
 - Extent
 - How to *impute* reasonable stand-in values
- Can be quite application specific/subjective!

More on this in future lectures!

Classification of postal codes [edit]		
Below is the list of postal codes in Taiwan.		
1 – Taipei [edit]		
Code	Division name	Chinese
Taipei City		
100	Zhongzheng District	中正區
103	Datong District	大同區
104	Zhongshan District	中山區
105	Songshan District	松山區
106	Da'an District	大安區
108	Wanhua District	萬華區
110	Xinyi District	信義區
111	Shilin District	士林區
112	Beitou District	北投區
114	Neihu District	內湖區
115	Nangang District	南港區
116	Wenshan District	文山區
2 – Keelung, New Taipei, Matsu, Yilan [edit]		
Code	Division name	Chinese
Keelung City		
200	Ren'ai District	仁愛區
201	Xinyi District	信義區
202	Zhongzheng District	中正區
203	Zhongshan District	中山區
204	Anle District	安樂區
205	Nuannuan District	暖暖區

Looking Back: Data Transformations

Modifying structure

- Splitting rows and columns (“splitrows”, “split”)
- Unnesting (“flatten”)

Hiding Things

- Rows (“keep”, “delete” a.k.a. “select”)
- Columns (“drop”, “aggregate”)

Adding Things

- “derive” a.k.a. “map”, “apply”

Changing Things

- “replace”/“set”

Text manipulation

- “extract”, “replace”, “set”

GroupBy/Aggregate arithmetic

- “aggregate”, a.k.a. “group by”, “reduce”

Multi-table Operations

- “join”, “lookup”
- “union”

SUGGESTIONS

keep
age_min
35
35
35
Affects all columns, 4481 rows

delete
age_min
35
35

Affects all columns, 5056 rows

set
age_min
35
30
40

Changes 1 column

derive
age_min
35
30
40

Affects 1 column, all rows

	age_min	column1
true	35	
false	30	
false	40	

Creates 1 column

Looking Back: Bigger Picture

What decisions did we make along the way?

- Data that got hidden
 - Choice of columns to drop
 - Values we cleaned: indicators, cleaned columns
 - Filtering of rows (vs. indicators)
- Data that got added
 - Other derived columns: calculations, splits, extractions
 - Joins, Unions
- Changes in granularity:
 - Coarser: grouping keys and aggregates
 - Finer: unnesting

How might the data wrangling influence the analyses we can do?

- Sometimes we won't figure this out until analysis begins
- And we loop back to wrangling!

Impediments to Collaboration

- Diversity of tools and PLs makes it hard to share
- Finding a script or computed result is harder than just writing the program from scratch!
 - Q: How could we fix this?
- View that much of the analysis work is “throw away”

Data Sources at Web Companies

- Examples from Facebook
 - Application databases
 - Web server logs
 - Event logs
 - API server logs
 - Ad server logs
 - Search server logs
 - Advertisement landing page content
 - Wikipedia
 - Images and video

Tabular Data

- What is a table?
 - A **table** is a collection of **rows** and **columns**
 - Each row has an **index**
 - Each column has a **name**
 - A **cell** is specified by an (index, name) pair
 - A cell may or may not have a **value**

Tabular Data

- Fortune 500

	A	B	C	D	E	F	G	H	I
1	rank	company	cik	ticker	sic	state_location	state_of_incorporation	revenues	profits
2	1	Wal-Mart Stores	104169	WMT	5331	AR	DE	421849	16389
3	2	Exxon Mobil	34088	XOM	2911	TX	NJ	354674	30460
4	3	Chevron	93410	CVX	2911	CA	DE	196337	19024
5	4	ConocoPhillips	1163165	COP	2911	TX	DE	184966	11358
6	5	Fannie Mae	310522	FNM	6111	DC	DC	153825	-14014
7	6	General Electric	40545	GE	3600	CT	NY	151628	11644
8	7	Berkshire Hathaway	1067983	BRKA	6331	NE	DE	136185	12967
9	8	General Motors	1467858	GM	3711	MI	MI	135592	6172
10	9	Bank of America Corp.	70858	BAC	6021	NC	DE	134194	-2238
11	10	Ford Motor	37996	F	3711	MI	DE	128954	6561
12	11	Hewlett-Packard	47217	HPQ	3570	CA	DE	126033	8761
13	12	AT&T	732717	T	4813	TX	DE	124629	19864
14	13	J.P. Morgan Chase & Co.	19617	JPM	6021	NY	DE	115475	17370
15	14	Citigroup	831001	C	6021	NY	DE	111055	10602
16	15	McKesson	927653	MCK	5122	CA	DE	108702	1263
17	16	Verizon Communications	732712	VZ	4813	NY	DE	106565	2549
18	17	American International Group	5272	AIG	6331	NY	DE	104417	7786
19	18	International Business Machines	51143	IBM	3570	NY	NY	99870	14833
20	19	Cardinal Health	721371	CAH	5122	OH	OH	98601.9	642.2
21	20	Freddie Mac	37785	FMC	2800	PA	DE	98368	-14025

Tabular Data

- Fortune 500

Fortune 500 with ticker and EDGAR ★

The screenshot shows a Google Sheets document titled "Fortune 500 with ticker and EDGAR". The spreadsheet contains a single sheet with four columns: C (cik), D (ticker), E (sic), and F (name). The data includes entries for companies like WMT, XOM, CVX, COP, FNM, GE, BRKA, GM, BAC, and F. A context menu is open, with the "Download as" option expanded, showing options for CSV, HTML, Text, Excel, OpenOffice, and PDF.

cik	ticker	sic	name
104169	WMT	5331	Wal-Mart Stores, Inc.
34088	XOM	2911	ExxonMobil
93410	CVX	2911	Chevron
1163165	COP	2911	Coca-Cola Company
310522	FNM	6111	Fannie Mae
40545	GE	3600	General Electric
1067983	BRKA	6331	Berkshire Hathaway
1467858	GM	3711	General Motors
70858	BAC	6021	Bank of America
37996	F	3711	First Data

Tabular Data (csv)

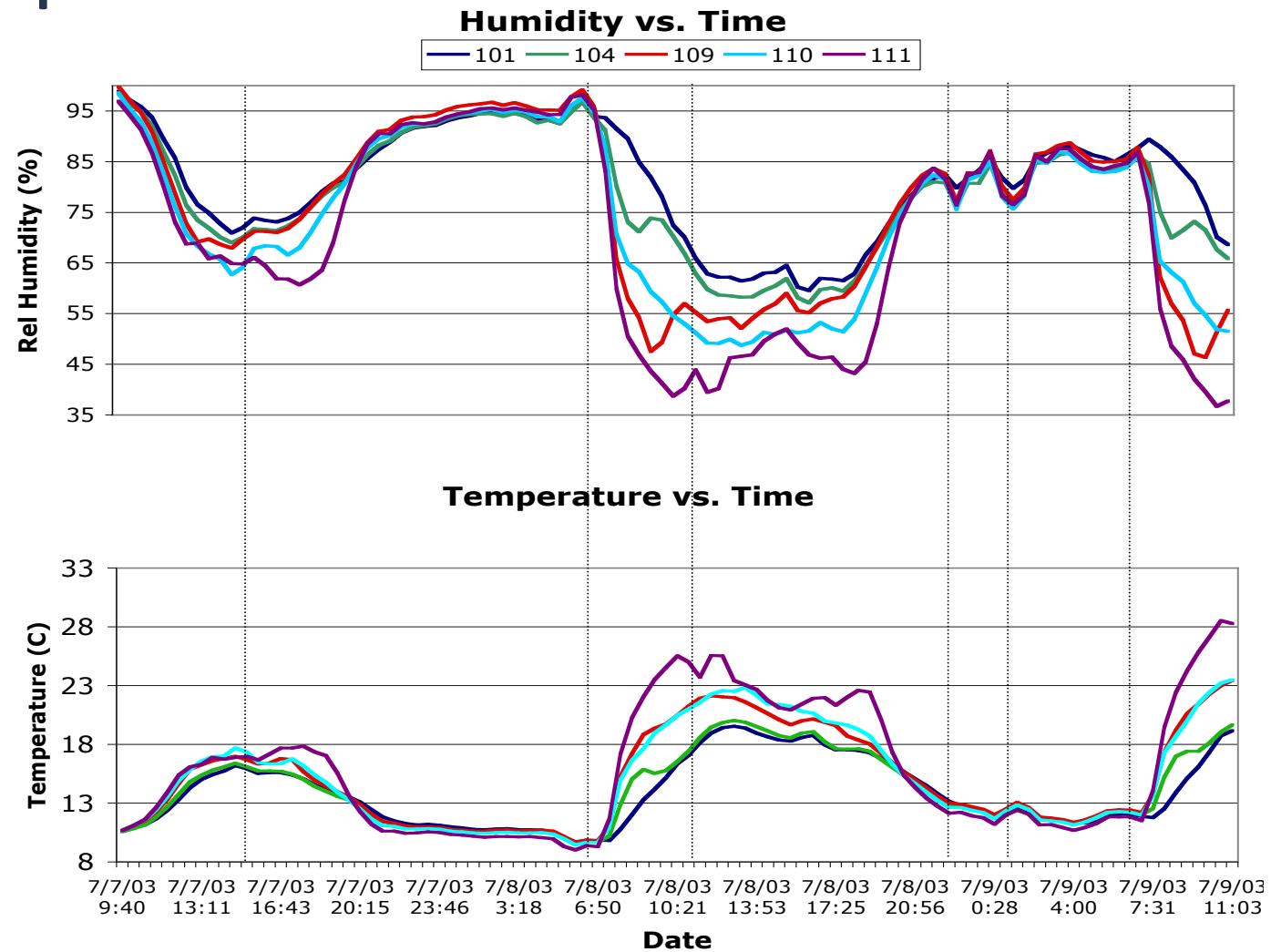
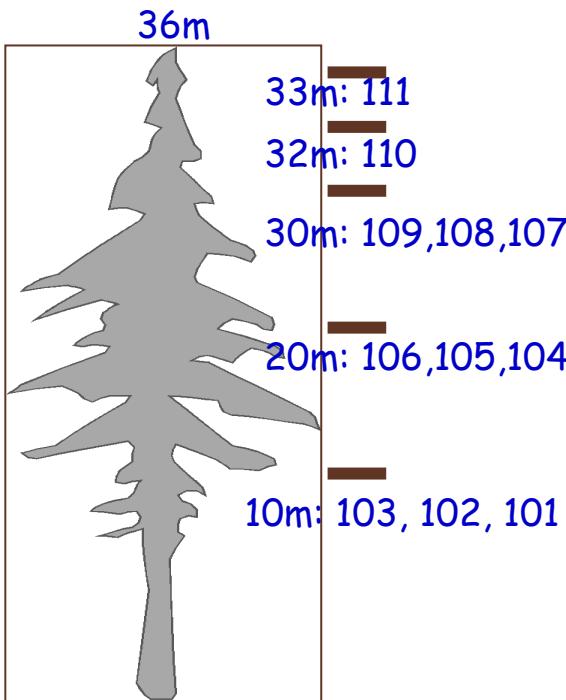
- Fortune 500

rank	company	cik	ticker	sic	state_location	state_of_incorporation	revenues	profits
1	Wal-Mart Stores	104169	WMT	5331	AR,DE	421849	16389	
2	Exxon Mobil	34088	XOM	2911	TX,NJ	354674	30460	
3	Chevron	93410	CVX	2911	CA,DE	196337	19024	
4	ConocoPhillips	1163165	COP	2911	TX,DE	184966	11358	
5	Fannie Mae	310522	FNM	6111	DC,DC	153825	-14014	
6	General Electric	40545	GE	3600	CT,NY	151628	11644	
7	Berkshire Hathaway	1067983	BRKA	6331	NE,DE	136185	12967	
8	General Motors	1467858	GM	3711	MI,MI	135592	6172	
9	Bank of America Corp.	70858	BAC	6021	NC,DE	134194	-2238	
10	Ford Motor	37996	F	3711	MI,DE	128954	6561	
11	Hewlett-Packard	47217	HPQ	3570	CA,DE	126033	8761	
12	AT&T	732717	T	4813	TX,DE	124629	19864	
13	J.P. Morgan Chase & Co.	19617	JPM	6021	NY,DE	115475	17370	
14	Citigroup	831001	C	6021	NY,DE	111055	10602	
15	McKesson	927653	MCK	5122	CA,DE	108702	1263	
16	Verizon Communications	732712	VZ	4813	NY,DE	106565	2549	
17	American International Group	5272	AIG	6331	NY,DE	104417	7786	
18	International Business Machines	51143	IBM	3570	NY,NY	99870	14833	
19	Cardinal Health	721371	CAH	5122	OH,OH	98601.9	642.2	
20	Freddie Mac	37785	FMC	2800	PA,DE	98368	-14025	
21	CVS Caremark	64803	CVS	5912	RI,DE	96413	3427	
22	UnitedHealth Group	731766	UNH	6324	MN,MN	94155	4634	
23	Wells Fargo	72971	WFC	6021	CA,DE	93249	12362	
24	Valero Energy	1035002	VLO	2911	TX,DE	86034	324	
25	Kroger	56873	KR	5411	OH,OH	82189.4	1116.3	
26	Procter & Gamble	80424	PG	2840	OH,OH	79689	12736	
27	AmerisourceBergen	1140859	ABC	5122	PA,DE	77954	636.7	
28	Costco Wholesale	909832	COST	5331	WA,WA	77946	1303	
29	Marathon Oil	101778	MRO	2911	TX,DE	68413	2568	
30	Home Depot	354950	HD	5211	GA,DE	67997	3338	

Protein Data Bank

HEADER APOPTOSIS 05-OCT-10 3IZA
TITLE STRUCTURE OF AN APOPTOSOME-PROCASPASE-9 CARD COMPLEX
COMPND MOL_ID: 1;
COMPND 2 MOLECULE: APOPTOTIC PROTEASE-ACTIVATING FACTOR 1;
COMPND 3 CHAIN: A, B, C, D, E, F, G;
COMPND 4 SYNONYM: APAF-1;
COMPND 5 ENGINEERED: YES
SOURCE MOL_ID: 1;
SOURCE 2 ORGANISM_SCIENTIFIC: HOMO SAPIENS;
SOURCE 3 ORGANISM_COMMON: HUMAN;
SOURCE 4 ORGANISM_TAXID: 9606;
SOURCE 5 GENE: APAF1, KIAA0413;
SOURCE 6 EXPRESSION_SYSTEM: SPODOPTERA FRUGIPERDA;
SOURCE 7 EXPRESSION_SYSTEM_TAXID: 7108;
SOURCE 8 EXPRESSION_SYSTEM_STRAIN: SF21;
SOURCE 9 EXPRESSION_SYSTEM_VECTOR_TYPE: INSECT VIRUS;
SOURCE 10 EXPRESSION_SYSTEM_PLASMID: PFASTBAC1
KEYWDS APOPTOSOME, APAF-1, PROCASPASE-9 CARD, APOPTOSIS
EXPDTA ELECTRON MICROSCOPY
AUTHOR S.YUAN,X.YU,M.TOPF,S.J.LUDTKE,X.WANG,C.W.AKEY
REVDAT 1 03-NOV-10 3IZA 0
SPRSDE 03-NOV-10 3IZA 3IYT
JRNL AUTH S.YUAN,X.YU,M.TOPF,S.J.LUDTKE,X.WANG,C.W.AKEY
JRNL TITL STRUCTURE OF AN APOPTOSOME-PROCASPASE-9 CARD COMPLEX
JRNL REF STRUCTURE V. 18 571 2010

Internet of Things: Example measurements



Tabular Data from Sensors



Challenges

- May be many missing fields (a particular sensor may not produce all types of output).
- Device may go offline for a while.
- Device may be damaged (permanently or intermittently).
- Timestamps usually critical but may not be accurate.
- Other meta-data (location, device ID) may have errors.

Log Files – Example Apache Web Log



Processes, usually daemons, create logs
e.g., `httpd`, `mysqld`, `syslogd`

- 66.249.65.107 - - [08/Oct/2007:04:54:20 -0400] "GET /support.html HTTP/1.1" 200 11179 "-" "Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
- 111.111.111.111 - - [08/Oct/2007:11:17:55 -0400] "GET / HTTP/1.1" 200 10801 "http://www.google.com/search?q=log+analyzer&ie=utf-8&oe=utf-8 &aq=t&rls=org.mozilla:en-US:official&client=firefox-a" "Mozilla/5.0 (Windows; U; Windows NT 5.2; en-US; rv:1.8.1.7) Gecko/20070914 Firefox/2.0.0.7"
- 111.111.111.111 - - [08/Oct/2007:11:17:55 -0400] "GET /style.css HTTP/1.1" 200 3225 "\"<http://www.loganalyzer.net/>" "Mozilla/5.0 (Windows; U; Windows NT 5.2; en-US; rv:1.8.1.7) Gecko/20070914 Firefox/2.0.0.7"

Syslog – A Standard for System Messages

- Developed by Eric Allman (at Berkeley) as part of the Sendmail project
- Standardized by the IETF in RFC 3164 and RFC 5424
- Listens on port 514 using UDP
- Puts data in /var/log/messages by default
- Enables rich analysis:

The Splunk logo features the word "splunk" in a large, bold, black sans-serif font. A registered trademark symbol (®) is positioned above the letter "k". To the right of the word, there is a grey chevron-like shape pointing to the right, consisting of three short diagonal lines.

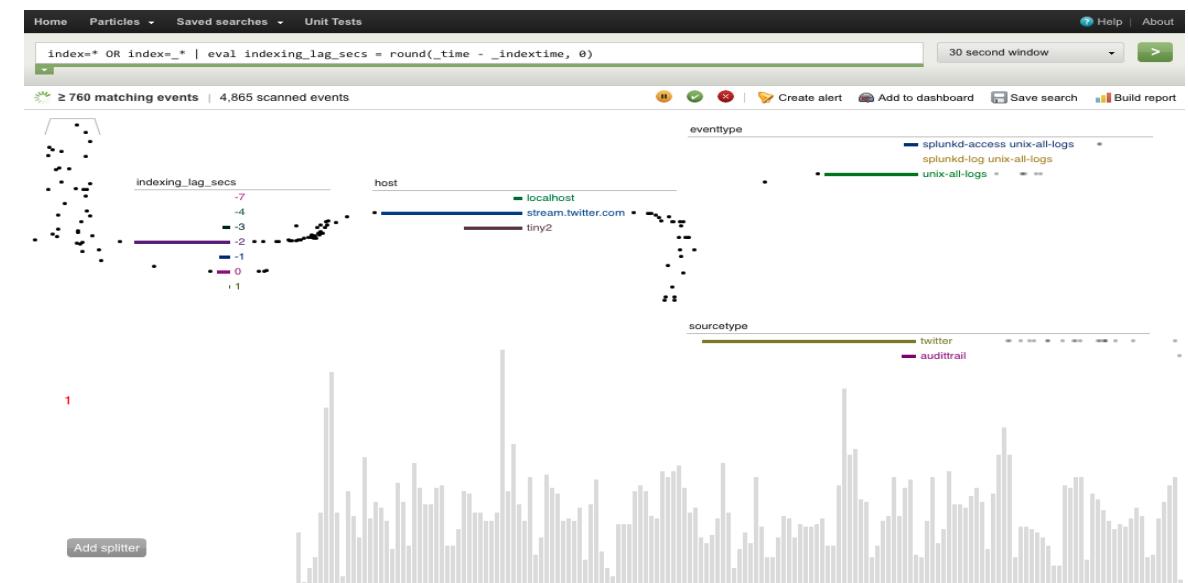
Syslog

dhcp-47-129:DataScienceF14> syslog -w 10

```
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMAccounting read:]: unexpected field ID 23 with type 8. Skipping.  
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMUser read:]: unexpected field ID 17 with type 12. Skipping.  
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMAuthenticationResult read:]: unexpected field ID 6 with type 11. Skipping.  
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMAuthenticationResult read:]: unexpected field ID 7 with type 11. Skipping.  
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMAccounting read:]: unexpected field ID 19 with type 8. Skipping.  
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMAccounting read:]: unexpected field ID 23 with type 8. Skipping.  
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMUser read:]: unexpected field ID 17 with type 12. Skipping.  
Feb 3 15:18:11 dhcp-47-129 Evernote[1140] <Warning>: -[EDAMSyncState read:]: unexpected field ID 5 with type 10. Skipping.  
Feb 3 15:18:49 dhcp-47-129 com.apple.mtmd[47] <Notice>: low priority thinning needed for volume Macintosh HD (/) with 18.9 <= 20.0 pct  
free space
```

“Splunking”

- Grab data from many machines
- Index it
- Check for unusual events:
 - Disk problems
 - Network congestion
 - Security attacks
- Monitor Resources
 - Network
 - Memory usage
 - Disk use, latency
 - Threads
- Dashboard for cloud administration.



Outline for this Evening

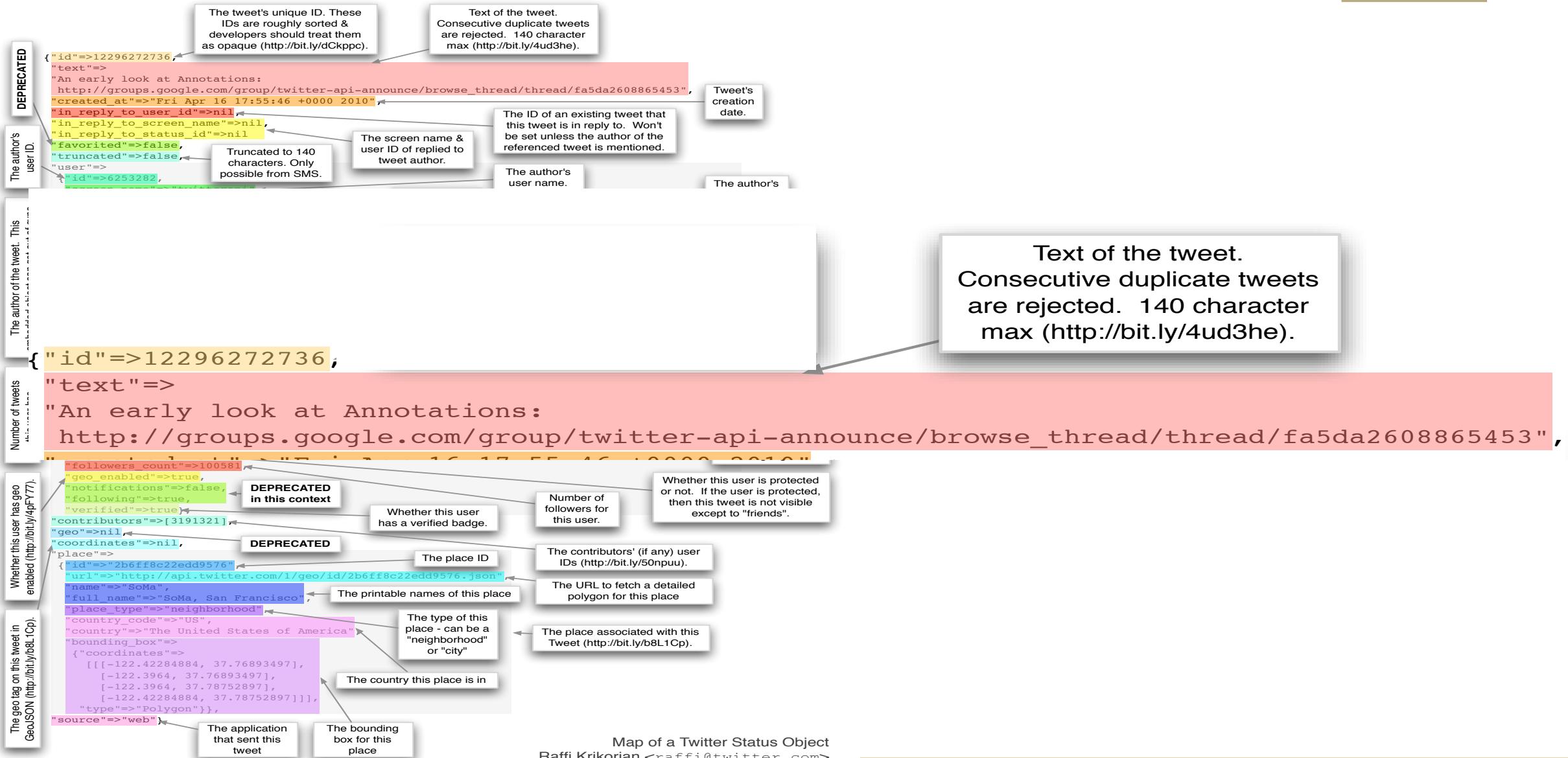
- Some Ideas from Kandel et al. Paper
- Data Types and Sources – “Document” data
- Data Preparation and Manipulation
- Files

Some Questions

- 1) How Many Characters are there in a Tweet?

- 2) How Many Bytes are there in a Tweet (msg)?

Tweet JSON Format



Tweet JSON Format



So how do we process the JSON from Twitter?

Stay tuned, but first some concepts from the XML world...

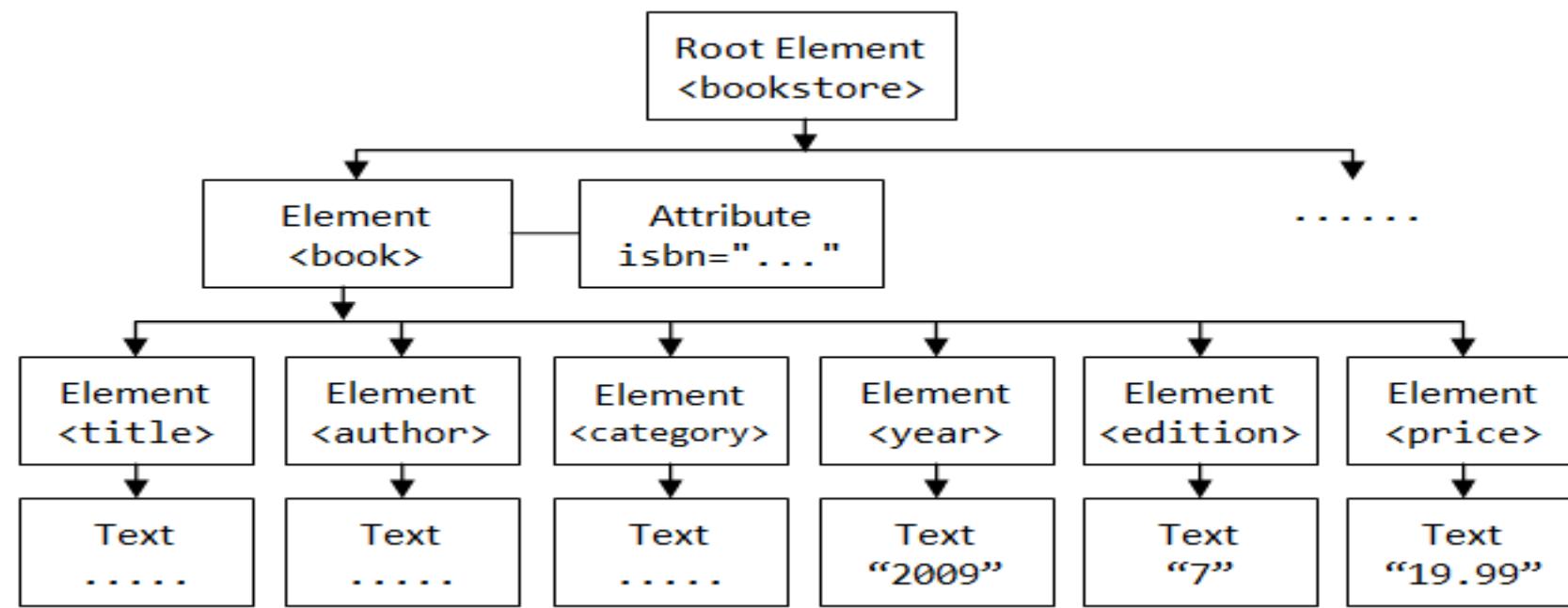
XML, DOM and SAX

```
<?xml version="1.0" encoding="UTF-8"?>
<!-- bookstore.xml -->
<bookstore>
  <book ISBN="0123456001">
    <title>Java For Dummies</title>
    <author>Tan Ah Teck</author>
    <category>Programming</category>
    <year>2009</year>
    <edition>7</edition>
    <price>19.99</price>
  </book>
```

XML, DOM and SAX

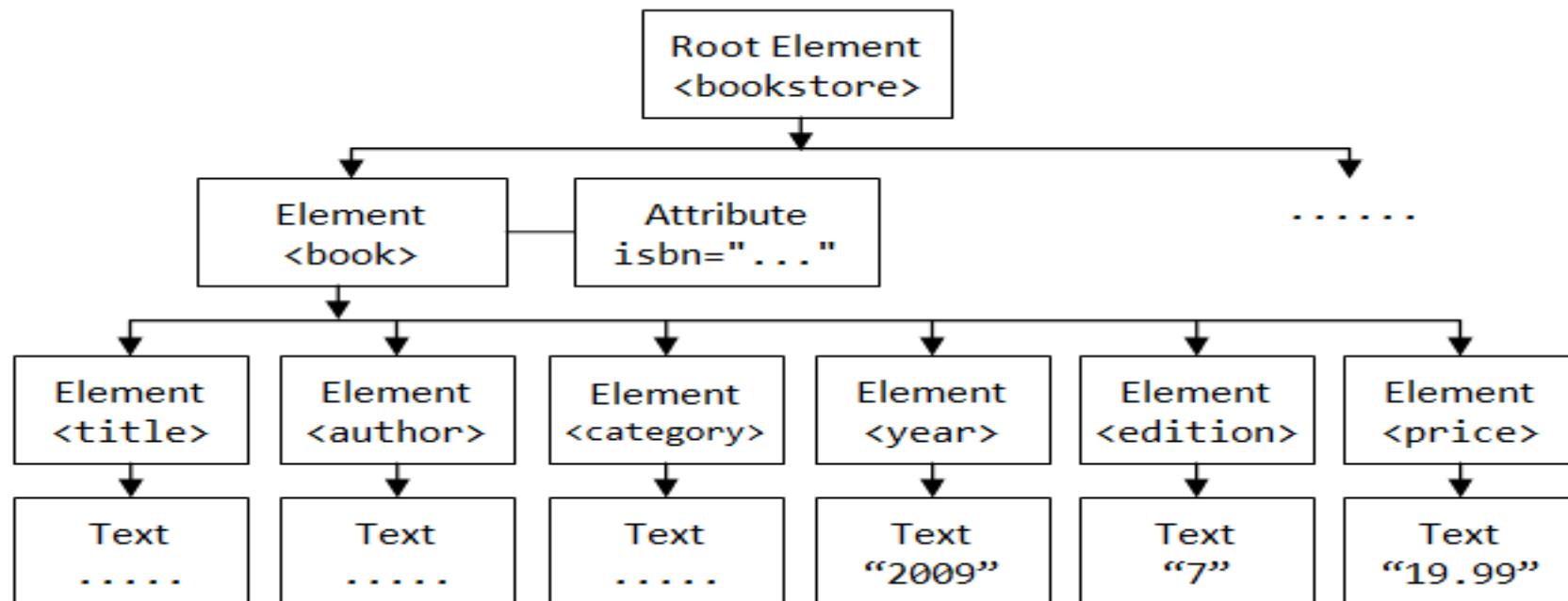
XML is a text format that encodes DOM (Document-Object Models) which is a data structure e.g. for Web pages.

The DOM is tree-structured:



XML, DOM and SAX

- The DOM is an easy object to work with: all the data in the object is accessible by links.
- The problem is that I might not care about most of the data, and I might not be able to fit the DOM for a large object in RAM.



SAX

SAX (Simple API for XML) is an alternative “lightweight” way to process XML.

A SAX parser generates a stream of events as it parses the XML file. The programmer registers handlers for each one.

It allows a programmer to handle only parts of the data structure.

SAX

```
<?xml version="1.0" encoding="UTF-8"?>          → Document Header
<!-- bookstore.xml -->           → Comment
<bookstore>           → Start-element "bookstore"
<book ISBN="0123456001">           → Start-element "book"
    <title>Java For Dummies</title>           → Start-element "title"
    <author>Tan Ah Teck</author>           → End-element "title"
    <category>Programming</category>
    <year>2009</year>
    <edition>7</edition>
    <price>19.99</price>
</book>           → End-element "book"
```

What about JSON?

Most JSON parsers construct the “DOM” directly.

But there are a few SAX-style parsers:

- Jackson
- JSON-simple

What about HTML?

- Common Crawl, **about 5 billion web pages**, between **0.2-0.5%** of Google's web crawl.
- 60 TB, hosted on Amazon S3, also available for download.
- Includes **link data**, **page rank**.
- In ARC (Internet Archive) File format.
- So there's plenty of data, and there are many crawlers for targeted exploration...
 - HTTrack, ...

HTML Tag Soup

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd">
<html xmlns="http://www.w3.org/1999/xhtml">
<head><!-- types/widgets/pages/common/page.tmpl home/index_v3.html generated by index_v3 on Wed 29 Feb 2012 11:04:41 PM PST -->
<title>San Francisco Bay Area &mdash; News, Sports, Business, Entertainment, Classifieds: SFGate</title>
<meta http-equiv="content-type" content="text/html; charset=iso-8859-1" />
<meta name="description" content="Find local news &amp; information, updated weather, traffic, classifieds, sports scores, real estate, jobs, cars, food &amp; wine, travel, entertainment, events and more on SFGate.com. Connect to the Bay Area community." />
<meta name="keywords" content="San Francisco, San Francisco Bay Area, news, local events, breaking news, world news, San Francisco Chronicle, SFGate" />
<meta property="fb:page_id" content="105702905593" />
<meta property="fb:admins" content="653226748,658759748" />
<!-- /widgets/sitewide/css/all/inc.html widgets/pages/common/post_write_mtime/css_inc.tmpl -->
<!-- generated by sitewidecss on Thu 16 Feb 2012 10:41:53 AM PST -->
<link rel="stylesheet" type="text/css" title="SFGate" media="all" href="http://imgs.sfgate.com/css1329417713/sitewide/css/sitewide.css" />
<!-- sitewide/css/all/inc.html end css_inc.tmpl -->
```

HTML Tools - Parsing

- “Beautiful Soup” <http://www.crummy.com/software/BeautifulSoup/>
a Python API for handling real HTML. DOM or SAX interfaces.
- “TagSoup”
<http://ccil.org/~cowan/XML/tagsoup/>
provides a Sax interface, i.e. a streaming parse, to Java applications. Can transform to a format you want using XSLT.
- Taggle, part of the Arabica toolset
<http://www.jezuk.co.uk/cgi-bin/view/arabica/code>
is a version of TagSoup written in C++. You probably want to use this if you have a lot of data.

Web Services

Most large web sites today actively discourage screen-scraping to get their content, and provide Web Service APIs instead.

This is the “right” way to get data from online sources.

Web Services

W3C definition:

a "Web service" as "a software system designed to support interoperable machine-to-machine interaction over a network".

Two kinds:

- XML-based RPC-style messages: WSDL and SOAP
- REST-style stateless interactions, URLs encode state

Can run over different transports, but usually HTTP

Examples

Twitter: REST API and streaming API with JSON content. Provides sampling, searching and filtering capabilities.

Amazon: has a “product advertising API” in XML with a WSDL spec. Includes product search, reviews etc.

Livejournal: RSS/Atom + custom XML/RPC. Search by keyword, topic, follow friend links.

Netflix: Javascript, Atom and REST interfaces.

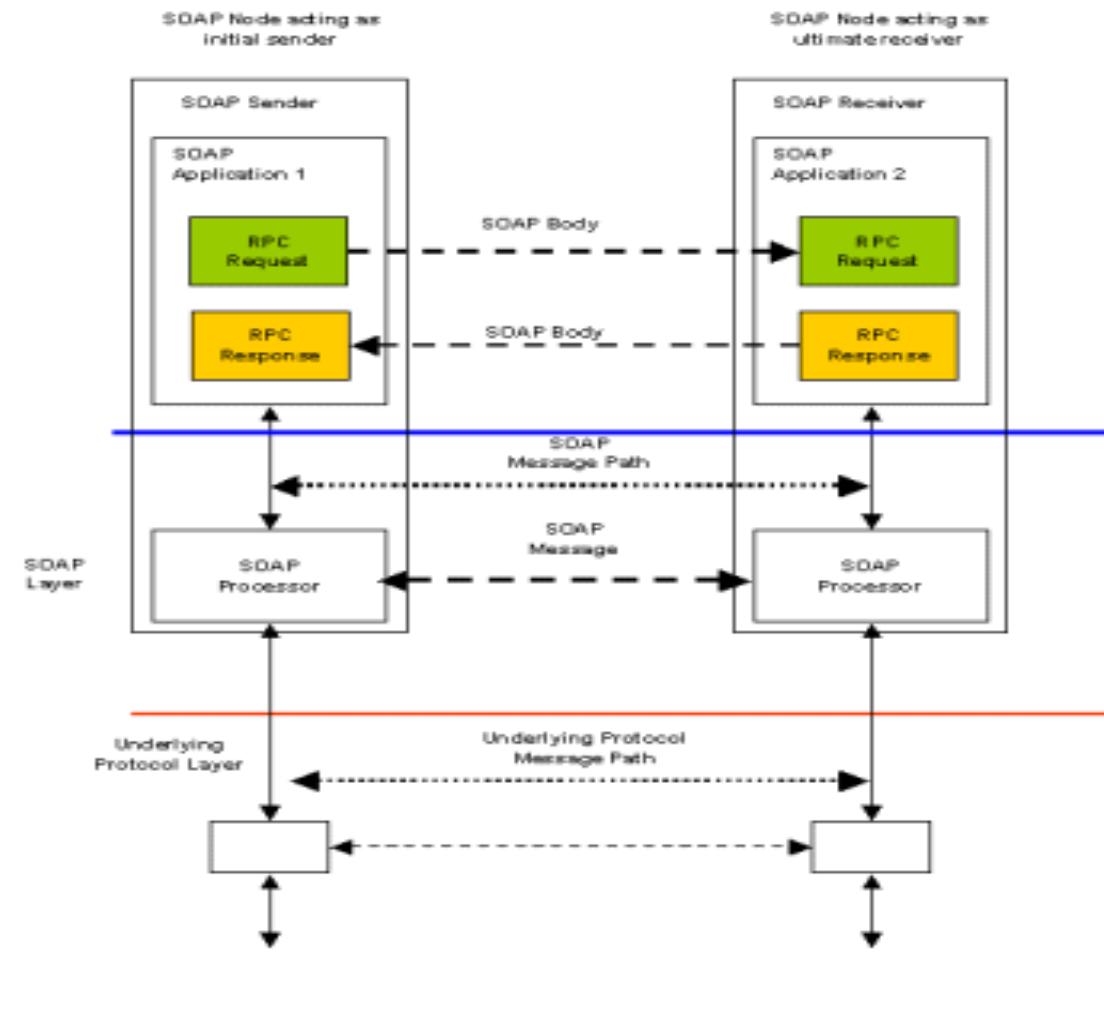
Ebay: Many APIs for searching, buying and posting. WSDL descriptions, client code in Java and .NET

Flickr: Comprehensive API set, free for non-commercial use. REST, XML-RPC, SOAP, with client code in many languages.

vBulletin: REST interface, most actions supported

Web Services

XML-RPC, requires a request-response cycle. Often longer “conversations.”



WSDL and SOAP

- Conceptually a Remote-Procedure-Call system, like CORBA, Java RMI etc.
- But RPC often has to pass through multiple layers of firewalls, causing many problems as security increased in the 1990s.
- Web services typically use HTTP as a transport, which runs almost anywhere, provides security and simple GET-POST methods. So HTTP-based RPC was a natural choice.
- SOAP uses XML messages, is human-readable, easy to process from any programming language, and relatively robust to version slippage.

SOAP is

SOAP covers the following four main areas:

- A **message format** for one-way communication describing how a message can be packed into an XML document.
- A **description** of how a SOAP message should be transported using HTTP (for Web-based interaction) or SMTP (for e-mail-based interaction).
- A **set of rules** that must be followed when processing a SOAP message and a simple classification of the entities involved in processing a SOAP message.
- A **set of conventions** on how to turn an RPC call into a SOAP message and back.

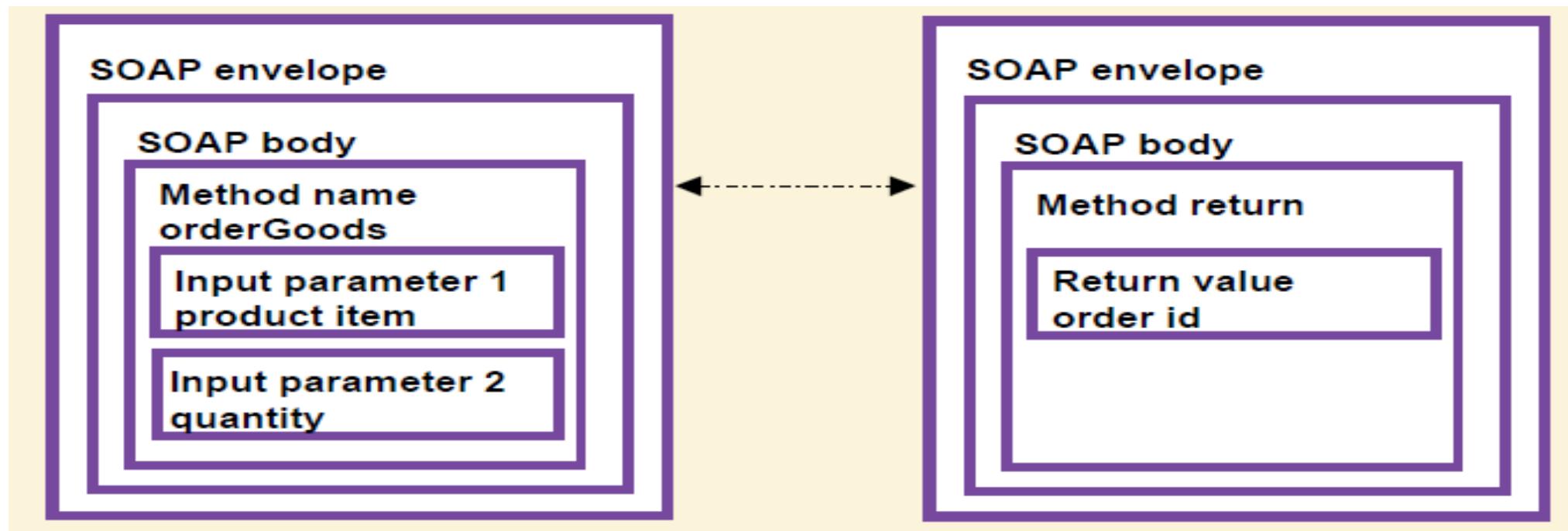
Soap Message

Typically an XML element containing header and body elements

```
<SOAP:Envelope xmlns:SOAP=
    "http://schemas.xmlsoap.org/soap/envelope/">
    <SOAP:Header>
        <!-- content of header goes here -->
    </SOAP:Header>
    <SOAP:Body>
        <!-- content of body goes here -->
    </SOAP:Body>
</SOAP:Envelope>
```

SOAP RPC

SOAP RPC messages typically encode arguments that are presented to the calling program as parameters and return values:



Soap RPC

```
POST /travelservice
SOAPAction: "http://www.acme-travel.com/flightinfo"
Content-Type: text/xml; charset="utf-8"
Content-Length: nnnn

<SOAP:Envelope xmlns:SOAP=
    "http://schemas.xmlsoap.org/soap/envelope/">
    <SOAP:Body>
        <m:GetFlightInfo
            xmlns:m="http://www.acme-travel.com/flightinfo"
            SOAP:encodingStyle=
                "http://schemas.xmlsoap.org/soap/encoding/"
            xmlns:xsd="http://www.w3.org/2001/XMLSchema"
            xmlns:xsi=
                "http://www.w3.org/2001/XMLSchema-instance">
            <airlineName xsi:type="xsd:string">UL
            </airlineName>
            <flightNumber xsi:type="xsd:int">506
            </flightNumber>
        </m:GetFlightInfo>
    </SOAP:Body>
</SOAP:Envelope>
```

Soap Response

```
HTTP/1.1 200 OK
Content-Type: text/xml; charset="utf-8"
Content-Length: nnnn

<SOAP:Envelope xmlns:SOAP=
    "http://schemas.xmlsoap.org/soap/envelope/">
    <SOAP:Body>
        <m:GetFlightInfoResponse
            xmlns:m="http://www.acme-travel.com/flightinfo"
            SOAP:encodingStyle=
                "http://schemas.xmlsoap.org/soap/encoding/"
            xmlns:xsd="http://www.w3.org/2001/XMLSchema"
            xmlns:xsi=
                "http://www.w3.org/2001/XMLSchema-instance">
            <flightInfo>
                <gate xsi:type="xsd:int">10</gate>
                <status xsi:type="xsd:string">ON TIME</status>
            </flightInfo>
        </m:GetFlightInfoResponse>
    </SOAP:Body>
</SOAP:Envelope>
```

REST

REpresentation State Transfer

Stateless Client/Server Protocol: Principles

1. Each message in the protocol contains all the information needed by the receiver to understand and/or process it. This constraint attempts to “keep things simple” and avoid needless complexity

2. Set of Uniquely Addressable Resources
 - “Everything is a Resource” in a RESTful system
 - Requires universal syntax for resource identification (e.g. URI)

REST

3. Set of Well-Defined Operations that can be applied to all resources

- In context of HTTP, the primary methods are
- POST, GET, PUT, DELETE
- these are similar (but not exactly) to the database notion of
- CRUD (Create, Read, Update, Delete)

4. The use of Hypermedia both for Application Information and State Transitions

- Resources are typically stored in a structured data format that supports hypermedia links, such as XHTML or XML

REST example

```
<user>
  <name>Jane</name>
  <gender>female</gender>
  <location href="http://www.example.org/us/ny/new_york">
    New York City, NY, USA</location>
</user>
```

This documentation is a representation used for the User resource

It might live at <http://www.example.org/users/jane/>

- If a user needs information about Jane, they GET this resource
- If they need to modify it, they GET it, modify it, and PUT it back
- The href to the Location resource allows savvy clients to gain access to its information with another simple GET request

Implication: Clients cannot be “thin”; need to understand resource formats

REST vs. RPC

In RPC systems, the design emphasis is on **verbs**

- What operations can I invoke on a system?
- `getUser()`, `addUser()`, `removeUser()`, `updateUser()`, `getLocation()`, `updateLocation()`, `listUsers()`, `listLocations()`, etc.

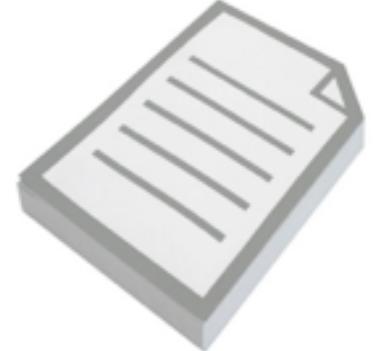
In REST systems, the design emphasis is on **nouns**

- User, Location
- In REST, you would define XML representations for these resources and then apply the standard methods to them

Outline for this Evening

- Some Ideas from Kandel et al. Paper
- Data Types and Sources
- Data Preparation and Manipulation
- Files

Files



- What is a file?
 - A **file** is a named sequence of **bytes**
 - Typically stored as a collection of pages (or blocks)
 - A **filesystem** is a collection of files organized within an hierarchical namespace
 - Responsible for laying out those bytes on physical media
 - Stores file metadata
 - Provides an API for interaction with files
 - Standard operations
 - `open()/close()`
 - `seek()`
 - `read()/write()`

Files

- Hierarchical namespace
 - / is known as the root of a filesystem
 - On Linux, the Filesystem Hierarchy Standard specifies which files live where
 - System executables in /usr/bin
 - Log files in /var/log
 - Permissions can be applied to all files beneath a directory
 - Files are not always arranged in a hierarchical namespace
 - Content-addressable storage (CAS)
 - Often used for large multimedia collections

File Formats

- Considerations for a file format
 - Data model: tabular, hierarchical, array
 - Physical layout
 - Field units and validation
 - Metadata: header, side file, specification, other?
 - Plain text or binary
 - Encoding: ASCII, UTF-8, other?
 - Delimiters and escaping
 - Compression, encryption, checksums?
 - Schema evolution

File Performance

Read/Write time (180 MB tabular file)

	Read Time (Text)	Write Time (Text)	Read Time (Binary)	Write Time (Binary)
Pandas (Python)	88 secs	107 secs	**	**
Scala/Java	12 secs	30 secs	0.5-2* secs	0.5-2* secs

** Pandas doesn't have a default binary file I/O library –
you can use Python, but performance depends on what you pick.

* 2 seconds is the time for sustainable read/write.
May be faster due to caching

File Performance - Compression

Read/Write time (180 MB tabular file, Scala/Java)

Binary File	Read Time	Write Time	File Size
Gzip level 6 (Java default)	1-2 secs	11 secs	35 MB
Gzip level 3	1-2 secs	4 secs	37 MB
Gzip level 1	1-2 secs	2 secs	38 MB
LZ4 fast	1-2 secs	1-2 secs	63 MB

Text File	Read Time	Write Time	File Size
Gzip level 6 (default)	1-2 secs	22 secs	32 MB
Gzip level 3	1-2 secs	14 secs	35 MB
Gzip level 1	1-2 secs	12 secs	39 MB
LZ4 fast	1-2 secs	11 secs	63 MB