

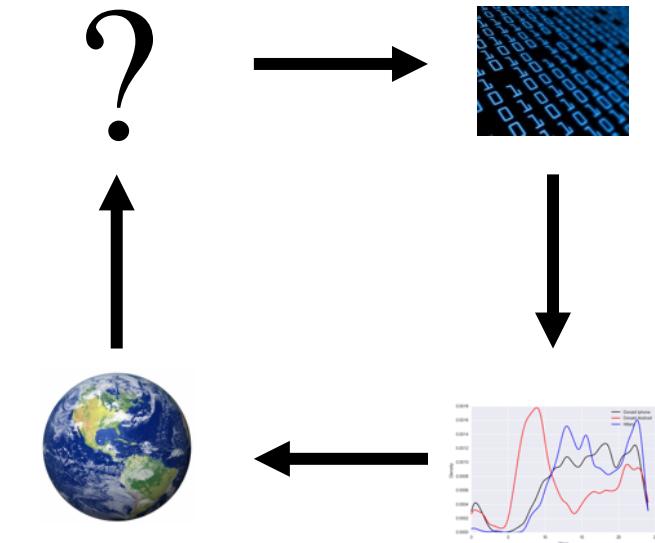
# Data Science 100

## *Lecture 4: Data Wrangling*

Slides by:

**Joe Hellerstein**

[hellerstein@berkeley.edu](mailto:hellerstein@berkeley.edu)



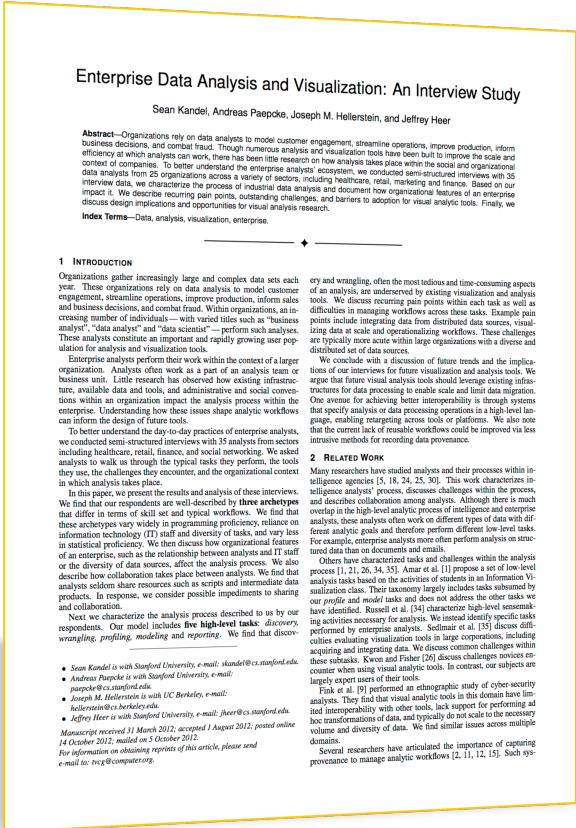
# So you want to be a data scientist...



# How will you spend your time?



# Enterprise Data Analysis and Visualization: An Interview Study



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.”



**Big Data  
Borat**

@BigDataBorat



Following

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.



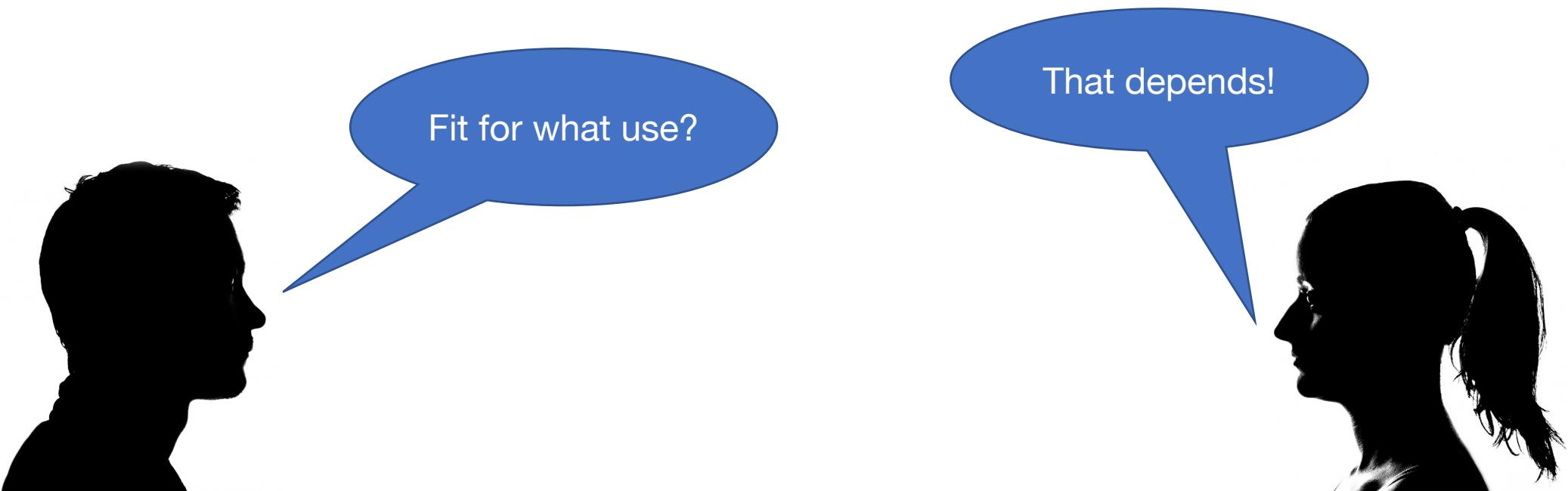
...

# 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



# Stages of Wrangling

## Raw: Data ingestion & discovery (“unboxing”)

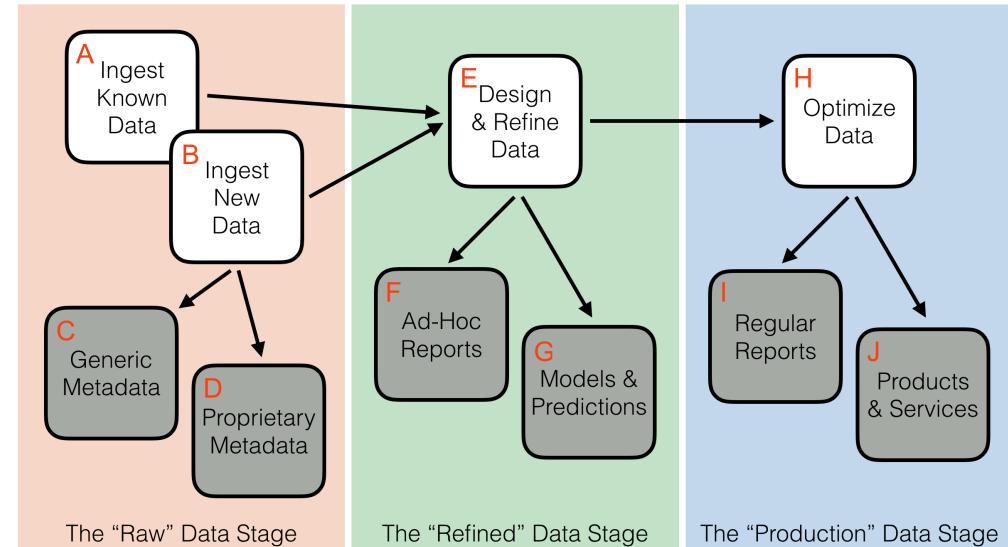
- What: Exploratory *ad hoc* analysis
- Who: The individual wrangler

## 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

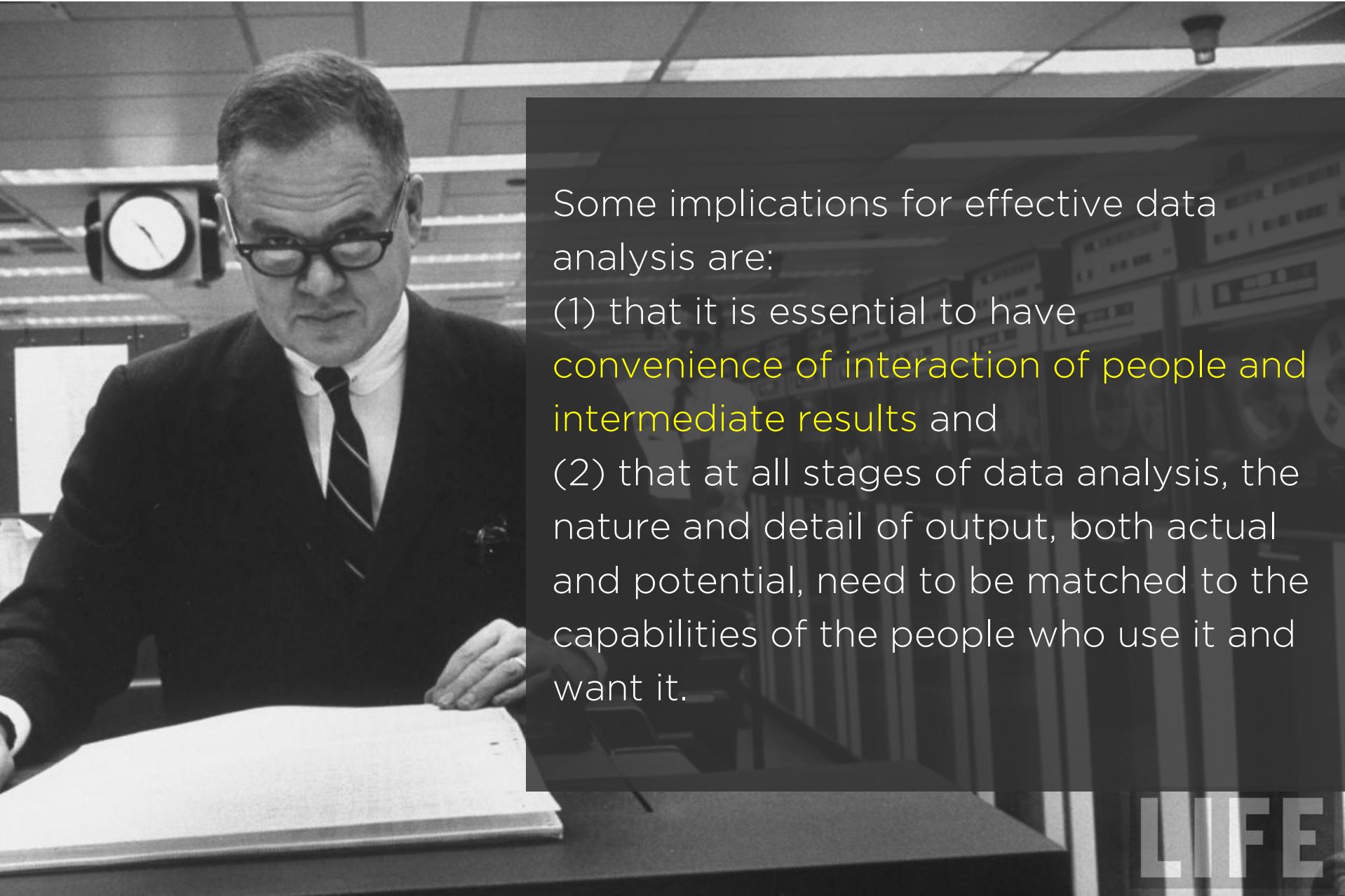
This is just a taste

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

More on this in future lectures!

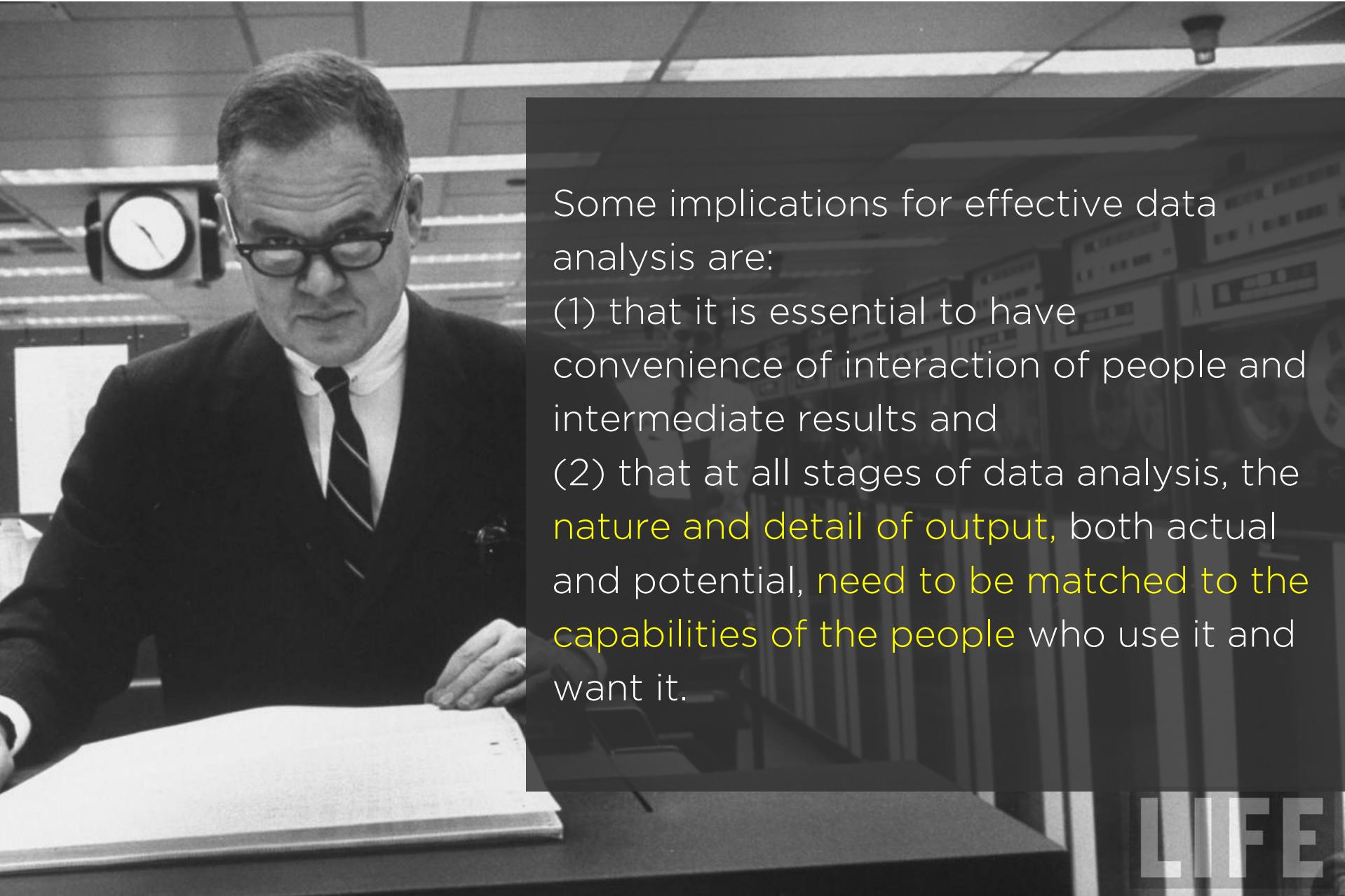


Data Analysis & Statistics, Tukey 1965



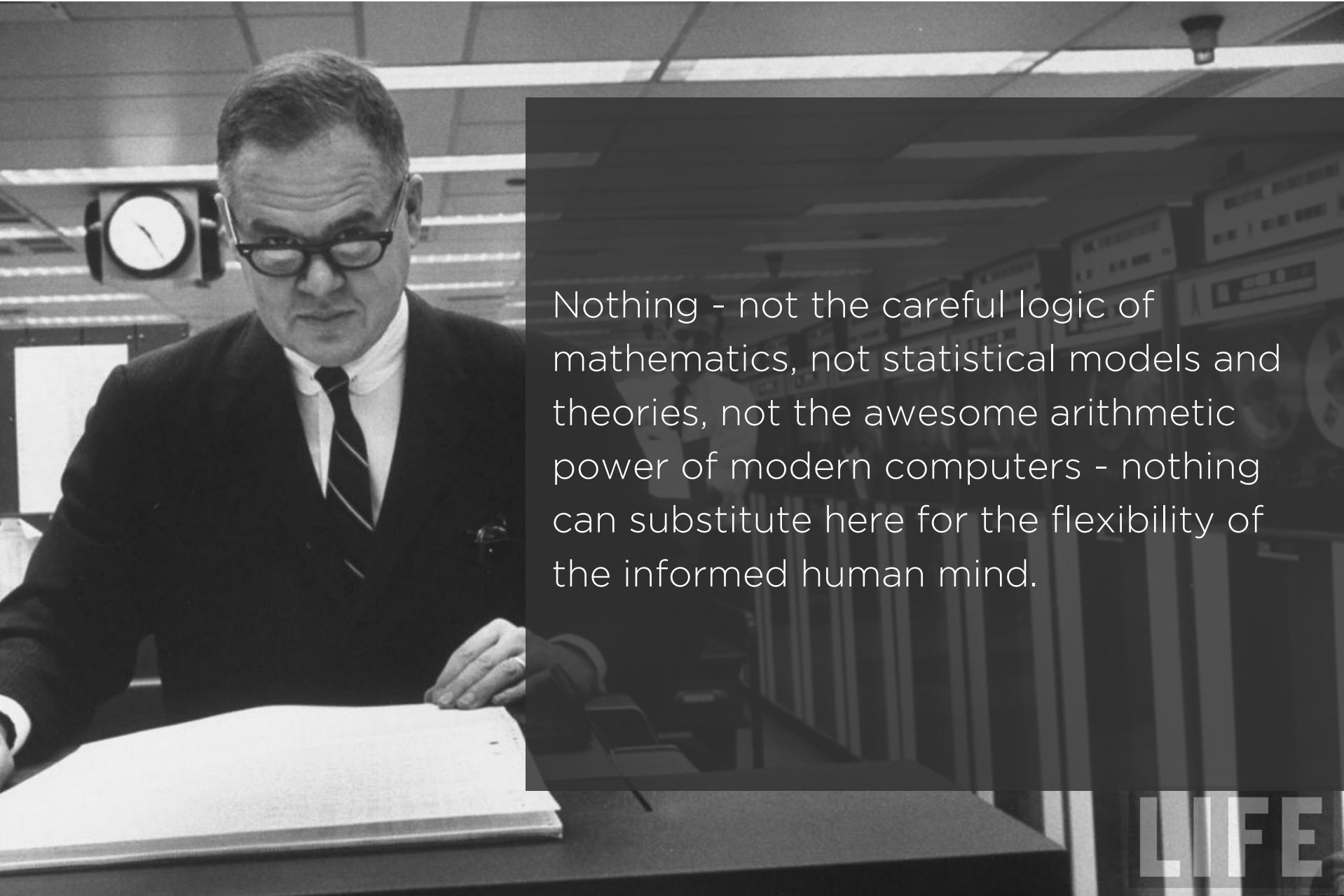
Some implications for effective data analysis are:

- (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.



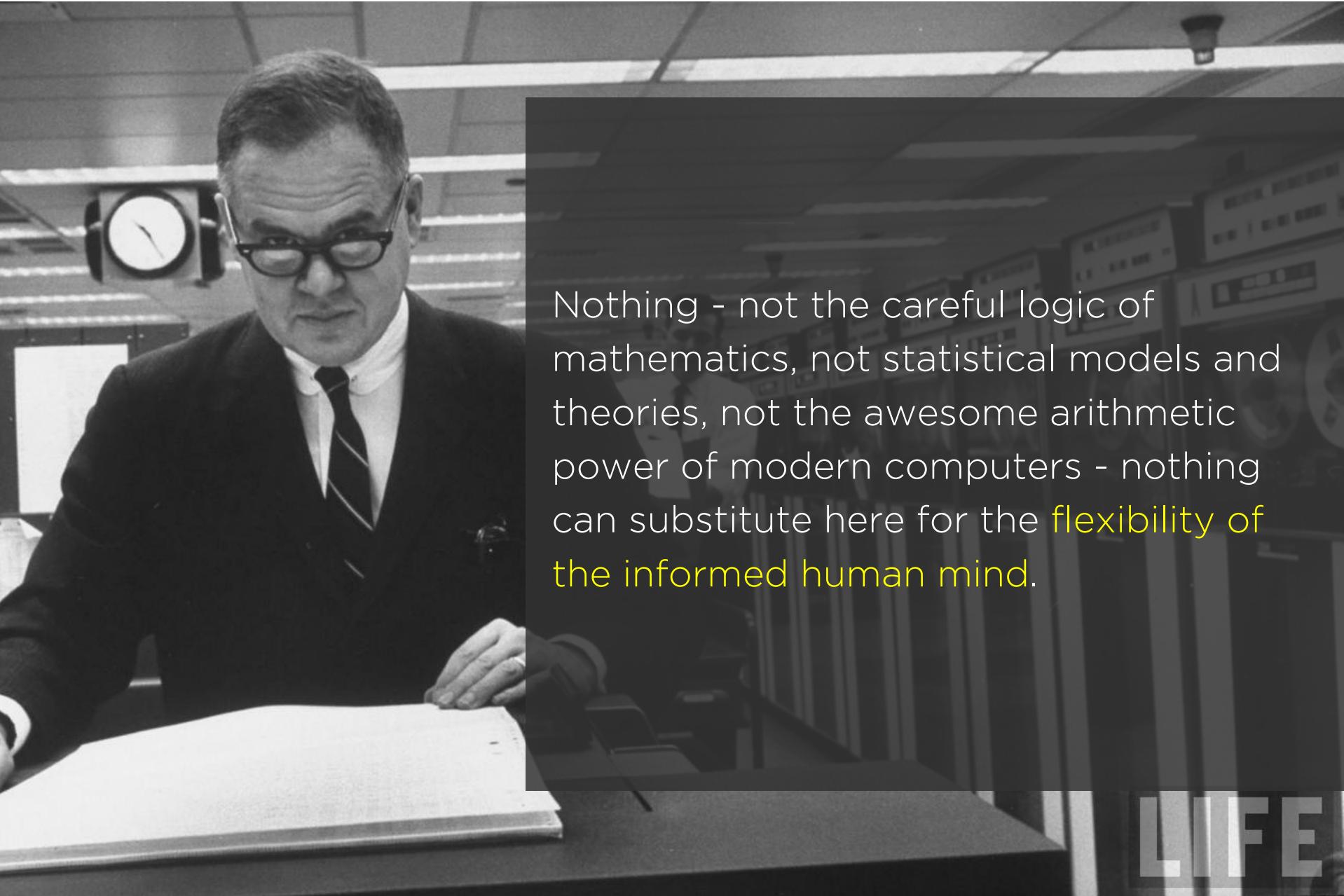
Some implications for effective data analysis are:

- (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.

LIFE



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.

LIFE

# 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*.

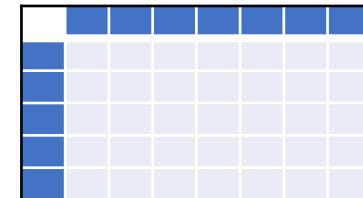
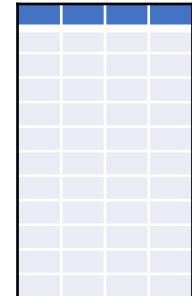
**(Data) Science is a human process.**

# 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

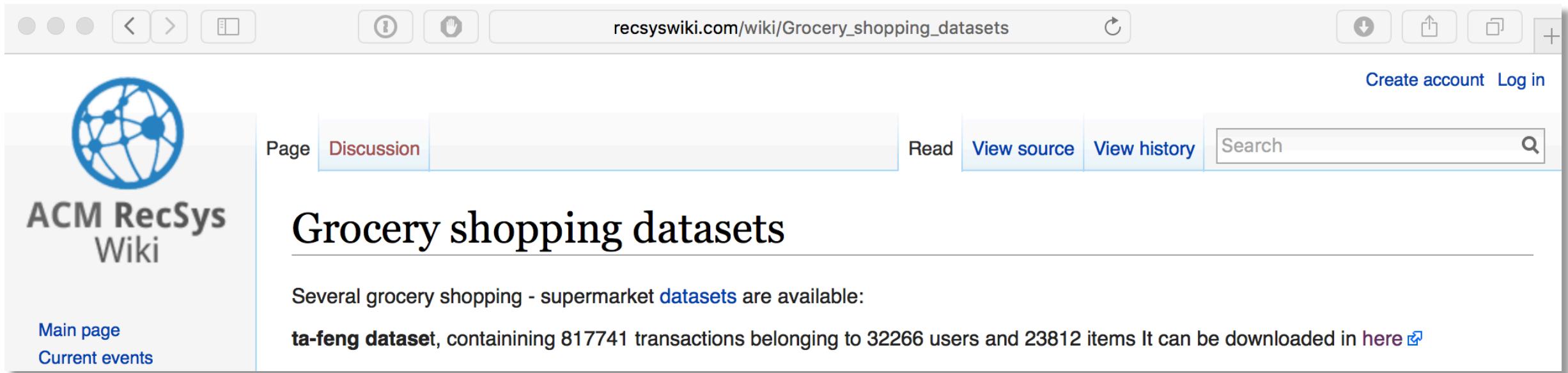


*There are non-rectangular structures as well, e.g. JSON, XML.*

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



# Without further ado, some data



The screenshot shows a web browser window with the following details:

- Address bar:** recsyswiki.com/wiki/Grocery\_shopping\_datasets
- Toolbar:** Standard browser controls for back, forward, search, and refresh.
- User Account:** Create account | Log in
- Page Header:** ACM RecSys Wiki
- Page Tabs:** Page (selected), Discussion, Read, View source, View history
- Search Bar:** Search
- Page Content:**
  - ## Grocery shopping datasets
  - Text: Several grocery shopping - supermarket [datasets](#) are available:
  - Text: **ta-feng dataset**, containing 817741 transactions belonging to 32266 users and 23812 items It can be downloaded in [here ↗](#)
- Sidebar:** Main page, Current events

Gently modified for pedagogical purposes [here](#).

Original dataset [here](#).

# Unboxing with UNIX Command Line

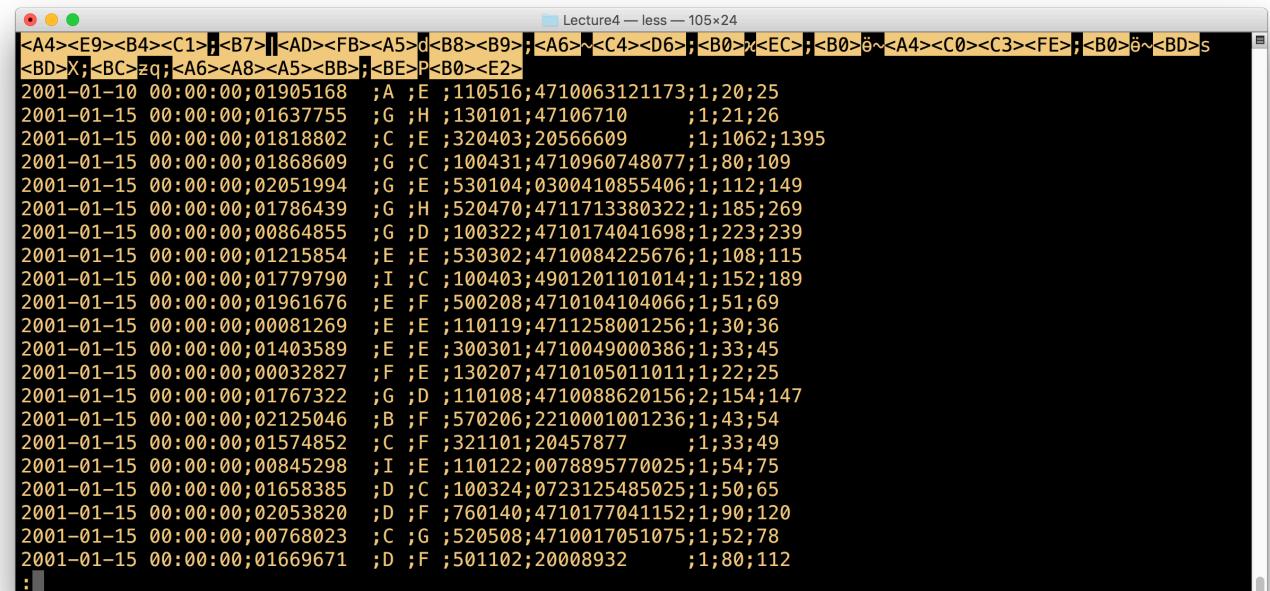
- File metadata
  - `ls -lh`
  - `file`
  - `WC`
- File (de)compression
  - `gunzip`, `zip`, `bzip`, etc.
- `stdout` and the pipe
- File content:
  - `cat`
  - `head`
  - `tail`
  - `less`
  - `<ctrl>-C`

```
[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
```



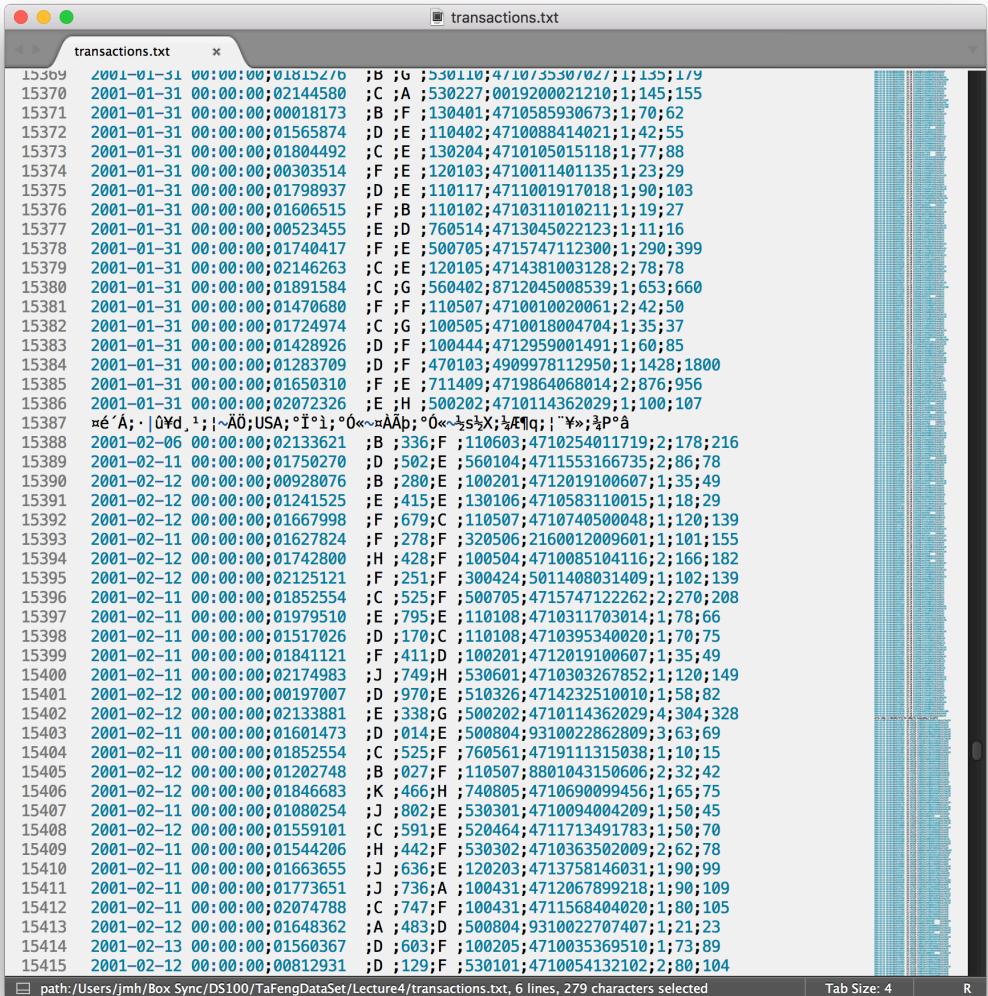
The screenshot shows a terminal window titled "Lecture4 — less — 105x24". The window displays a series of records, each consisting of a timestamp, a sequence of characters, and a semicolon-separated list of values. The timestamp is in the format "2001-01-15 00:00:00". The sequence of characters is mostly composed of lowercase letters (A through Z). The values are separated by semicolons and include various numerical and textual data. The terminal window has a dark background with light-colored text.

```
<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;252001-01-15 00:00:00;01637755 ;G ;H ;130101;47106710 ;1;21;262001-01-15 00:00:00;01818802 ;C ;E ;320403;20566609 ;1;1062;13952001-01-15 00:00:00;01868609 ;G ;C ;100431;4710960748077;1;80;1092001-01-15 00:00:00;02051994 ;G ;E ;530104;0300410855406;1;112;1492001-01-15 00:00:00;01786439 ;G ;H ;520470;4711713380322;1;185;2692001-01-15 00:00:00;00864855 ;G ;D ;100322;4710174041698;1;223;2392001-01-15 00:00:00;01215854 ;E ;E ;530302;4710084225676;1;108;1152001-01-15 00:00:00;01779790 ;I ;C ;100403;4901201101014;1;152;1892001-01-15 00:00:00;01961676 ;E ;F ;500208;4710104104066;1;51;692001-01-15 00:00:00;00081269 ;E ;E ;110119;4711258001256;1;30;362001-01-15 00:00:00;01403589 ;E ;E ;300301;4710049000386;1;33;452001-01-15 00:00:00;00032827 ;F ;E ;130207;4710105011011;1;22;252001-01-15 00:00:00;01767322 ;G ;D ;110108;4710088620156;2;154;1472001-01-15 00:00:00;02125046 ;B ;F ;570206;2210001001236;1;43;542001-01-15 00:00:00;01574852 ;C ;F ;321101;20457877 ;1;33;492001-01-15 00:00:00;00845298 ;I ;E ;110122;0078895770025;1;54;752001-01-15 00:00:00;01658385 ;D ;C ;100324;0723125485025;1;50;652001-01-15 00:00:00;02053820 ;D ;F ;760140;4710177041152;1;90;1202001-01-15 00:00:00;00768023 ;C ;G ;520508;4710017051075;1;52;782001-01-15 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!



# 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



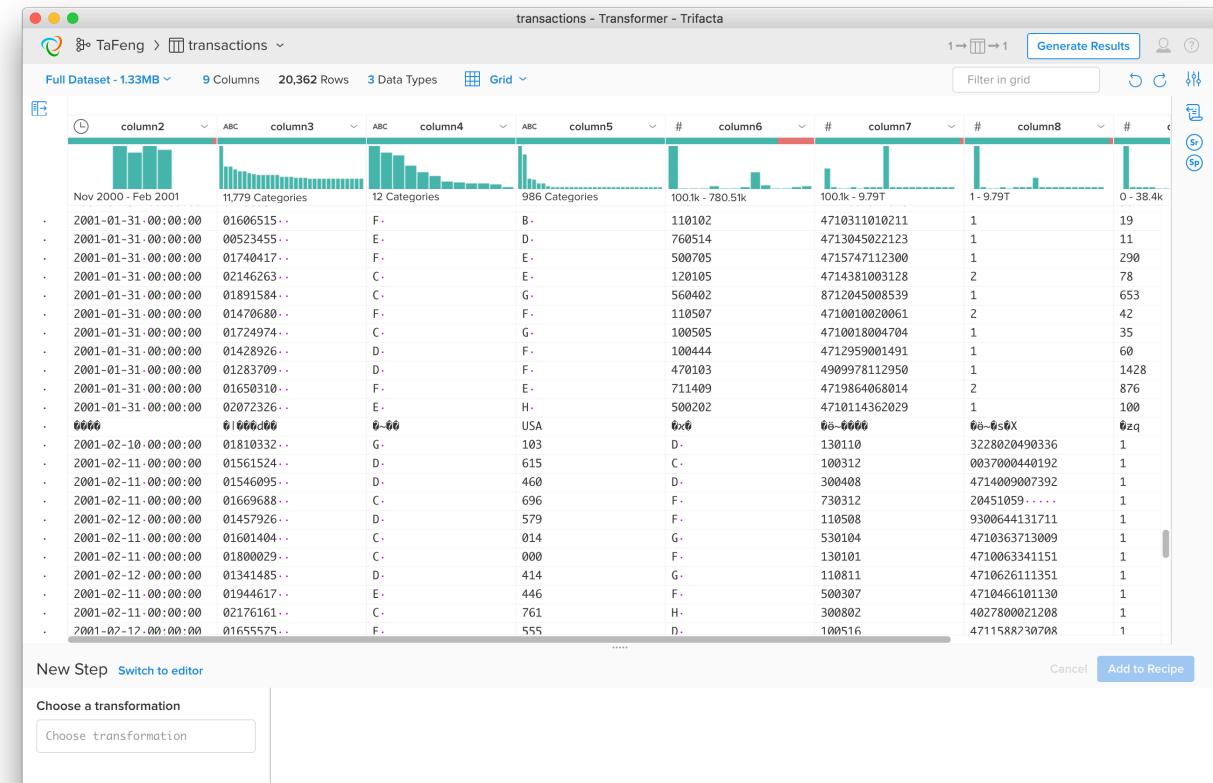
- Visual Profiling + Transformation
  - Iterative and ongoing
- Predictive Interaction
- Transformation language: Wrangle
  - Superset of relational algebra
  - Wrangle (and Pandas) both related to SQL
  - Generates a “recipe” (script)
    - You can run it on data to generate new data
- Scale & Interoperability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



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.

# 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

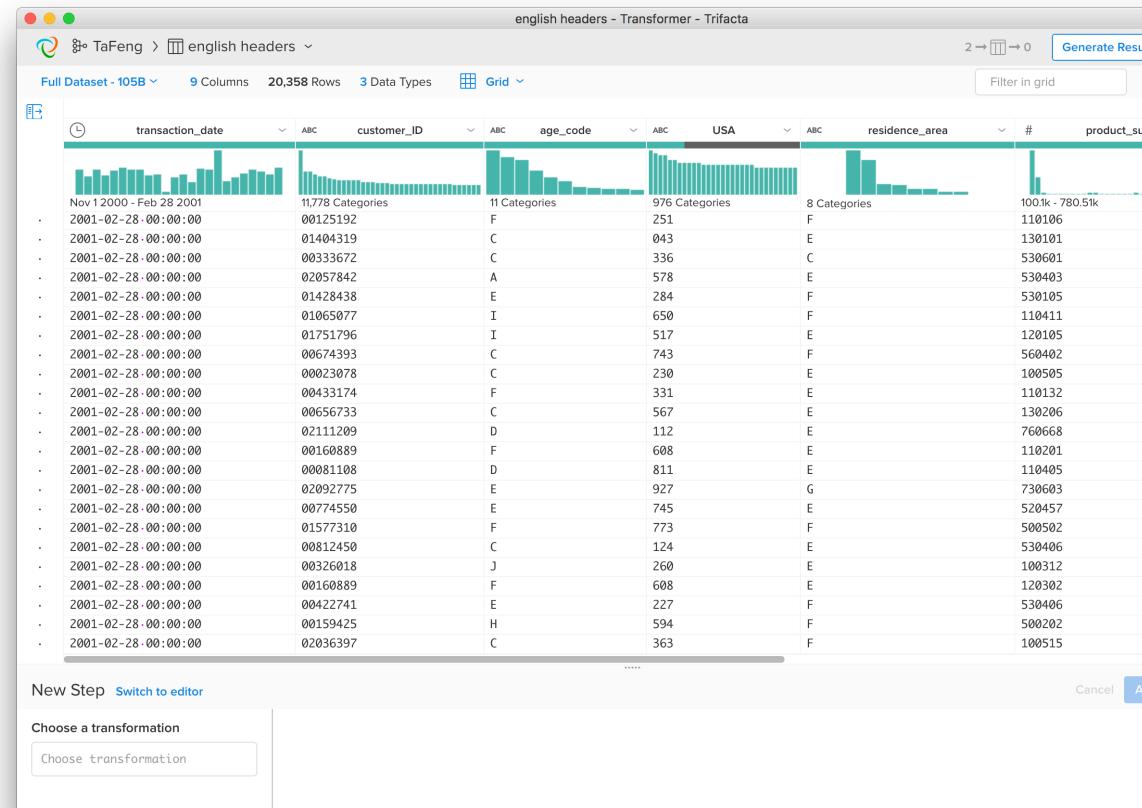
## Numbers and Strings

*We'll spend time here today*

## Common special-case “primitives”

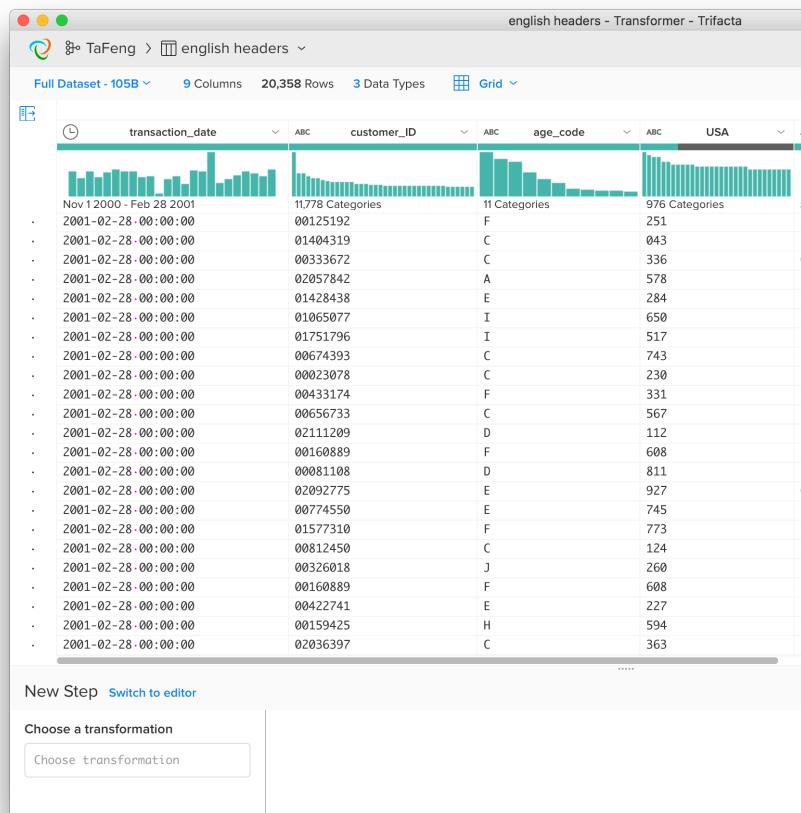
- Date/Time
- Geolocations

*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?

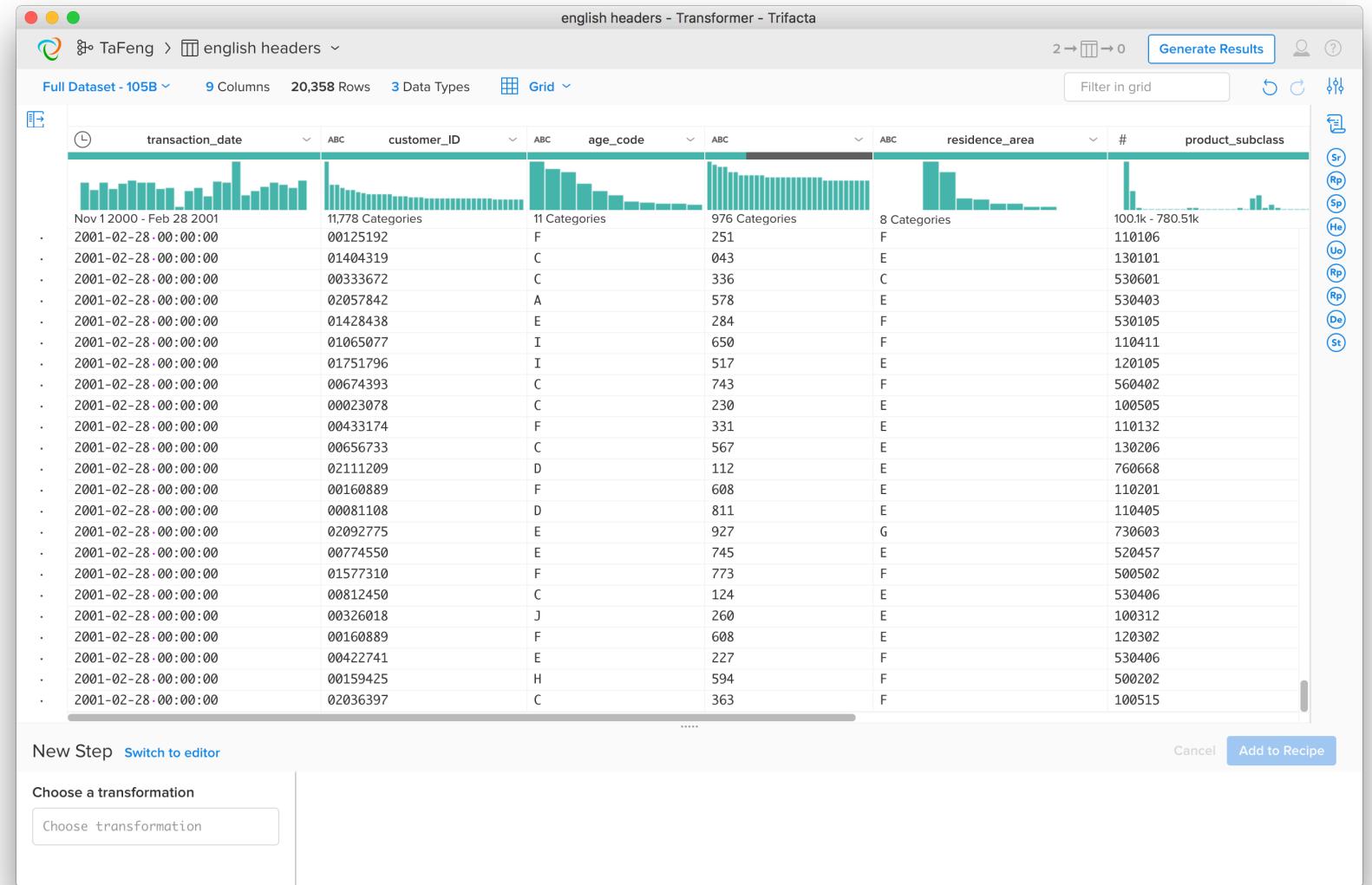
The screenshot shows a data visualization tool with the following details:

- Sample 1 - First 500kB:** 25 Columns, 140 Rows, 4 Data Types, Grid.
- Metrics:** 140 Categories, 0.85%.
- Columns:** id, favorite\_count, source.
- Data Preview:** A list of 55 tweets, each with an ID, a favorite count, and a source URL. The sources include various Twitter clients and web browsers.

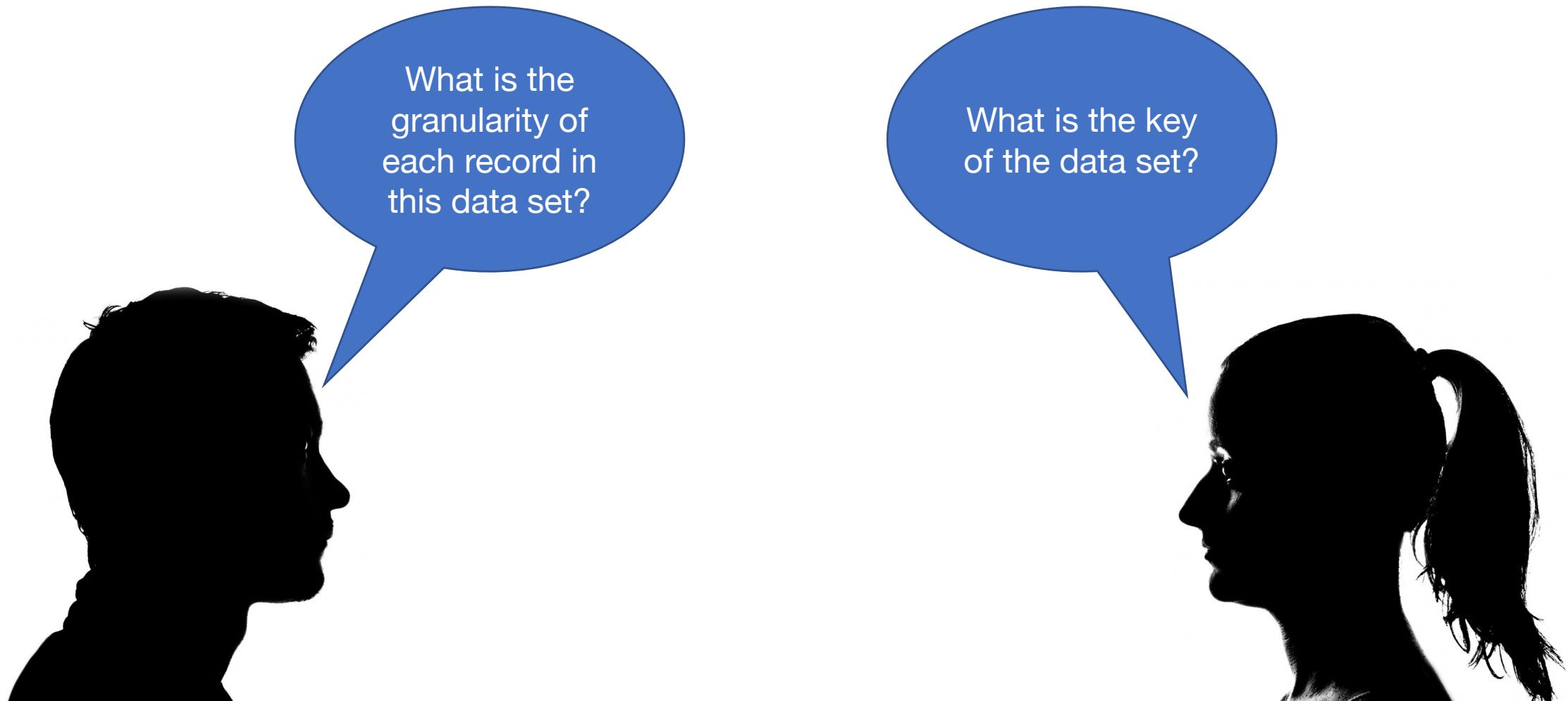
id	favorite_count	source
1	763959742755434500	7000
2	7639497233318700	7542
3	76394972331845400	5333
4	763857903299537000	12575
5	763748672012735600	25779
6	763558599339191000	20478
7	763573272954388000	31135
8	76351964865310000	41672
9	76351964441000000	18460
10	76351649065523000	19863
11	76348180225450000	14104
12	76347408905004000	19188
13	76347491183167000	19475
14	76345452693120000	16823
15	76345452693120000	20787
16	76339630123110000	20251
17	763391459118314000	27981
18	763385288295054000	17173
19	763385288295054000	36468
20	763385288295054000	45699
21	763167670174232000	18125
22	763167345405784000	11878
23	7631195713099981000	23645
24	7631195713099981000	0
25	7631195713099981000	0
26	7631195713099981000	0
27	763108383899777000	0
28	7629823608530505000	17882
29	762982142783764000	13658
30	762982142783764000	44682
31	7629714890000532000	11296
32	7628421269000152000	25518
33	7628157554235187000	0
34	7627906713260500000	17211
35	7627818262649605000	33203
36	7627771664762950000	25247
37	7627771664762950000	16379
38	7627764406810952000	15253
39	762775525765143000	13452
40	76277432138036794000	13116
41	76277432138036794000	11893
42	76277432138036794000	12983
43	762743923697722000	0
44	76269882571988000	21419
45	76264159537139190000	13946
46	7625391900112130000	19351
47	7625371321550500000	26460
48	7624008690581150000	41271
49	762399109407065000	0
50	7622845533416750000	33375
51	762210918721107000	20572
52	762210918721107000	21188
53	762184411795611000	8596
54	762010426182296000	25958
55		

# What does each TaFeng record represent?

Theories?  
Justification?



# Assessing Granularity



# Key (“primary key”)



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)`

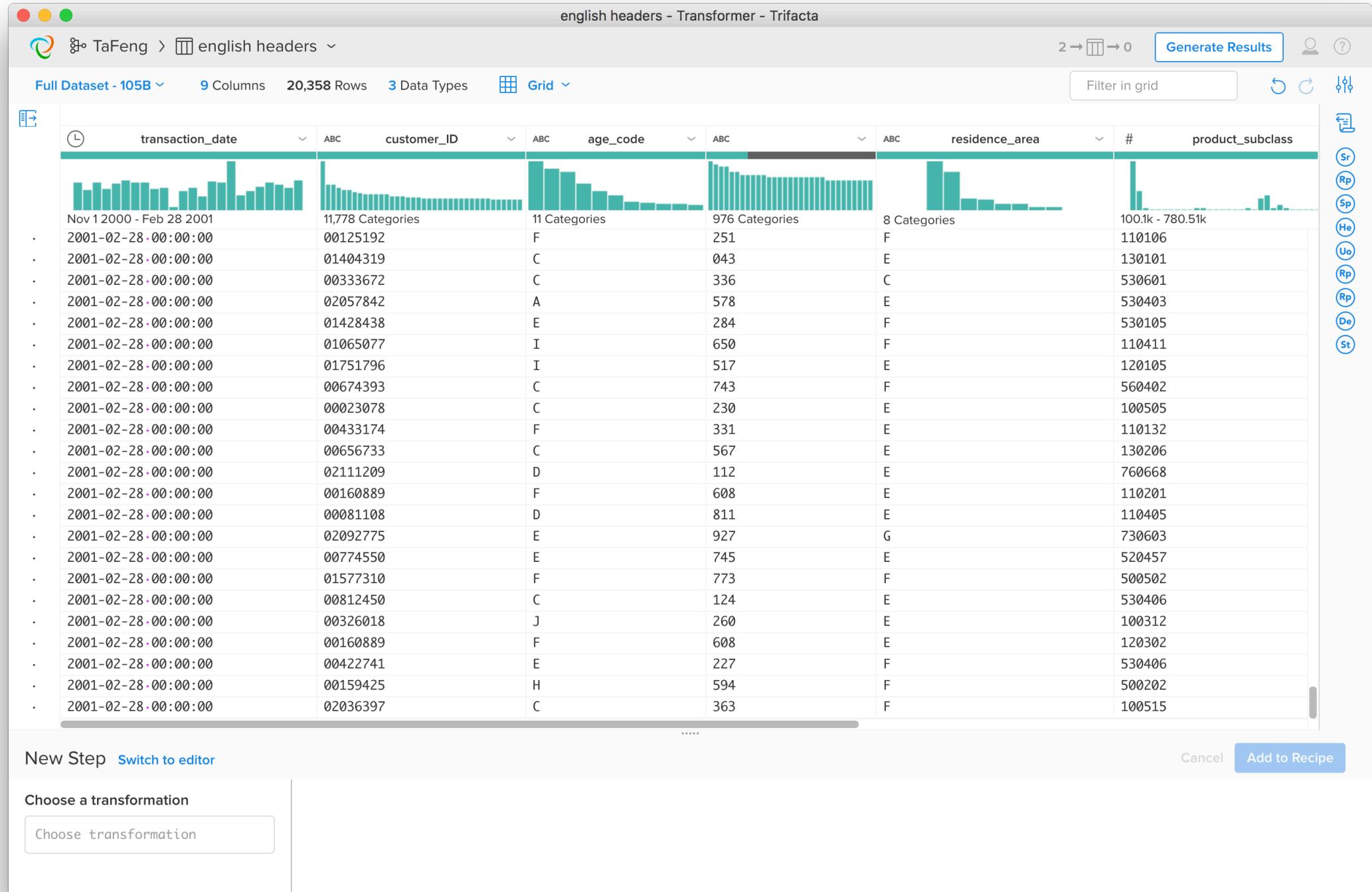
Each value occurs **at most once** in the key column(s).

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



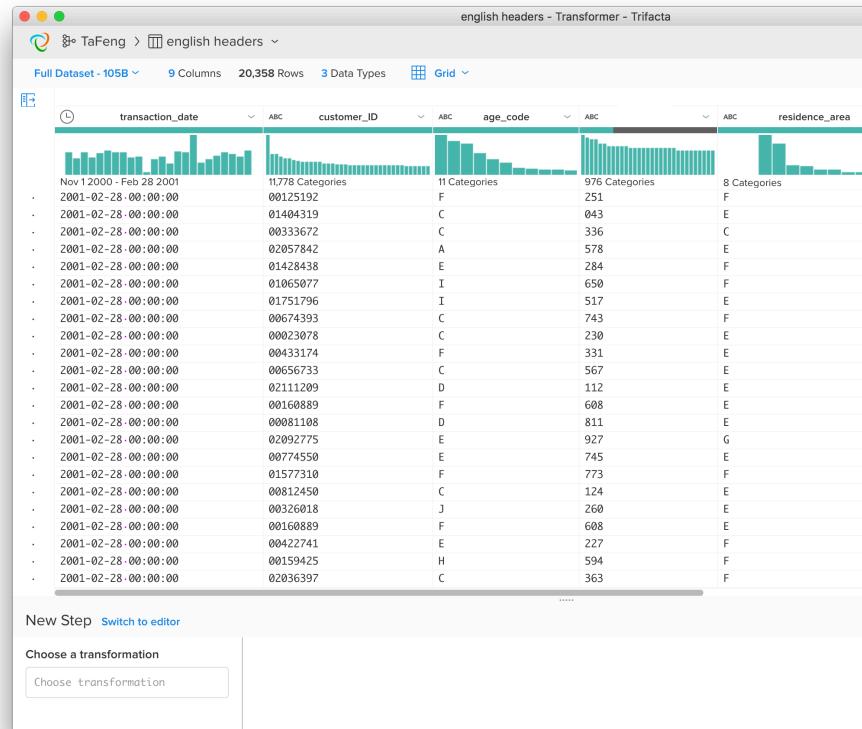
# Joining tables with keys

- 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

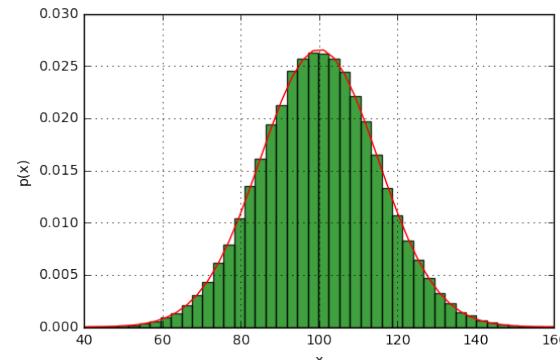
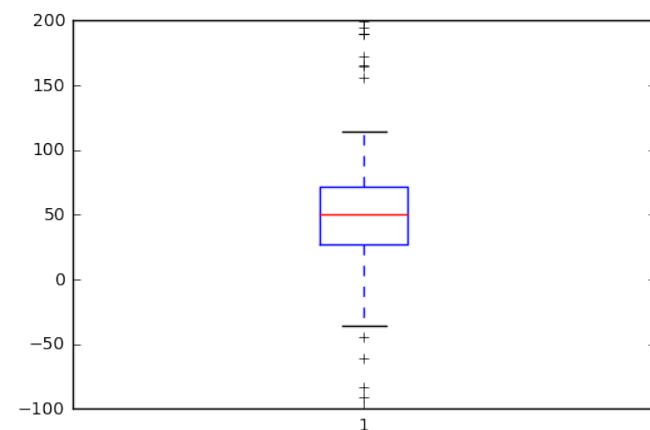
What is an “outlier”?

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

More on this in future lectures!

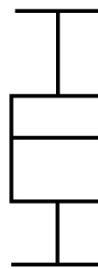
Distribution-based definition

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



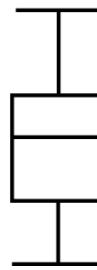
# What to do with Outliers?

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	Count
20002000	8

**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**

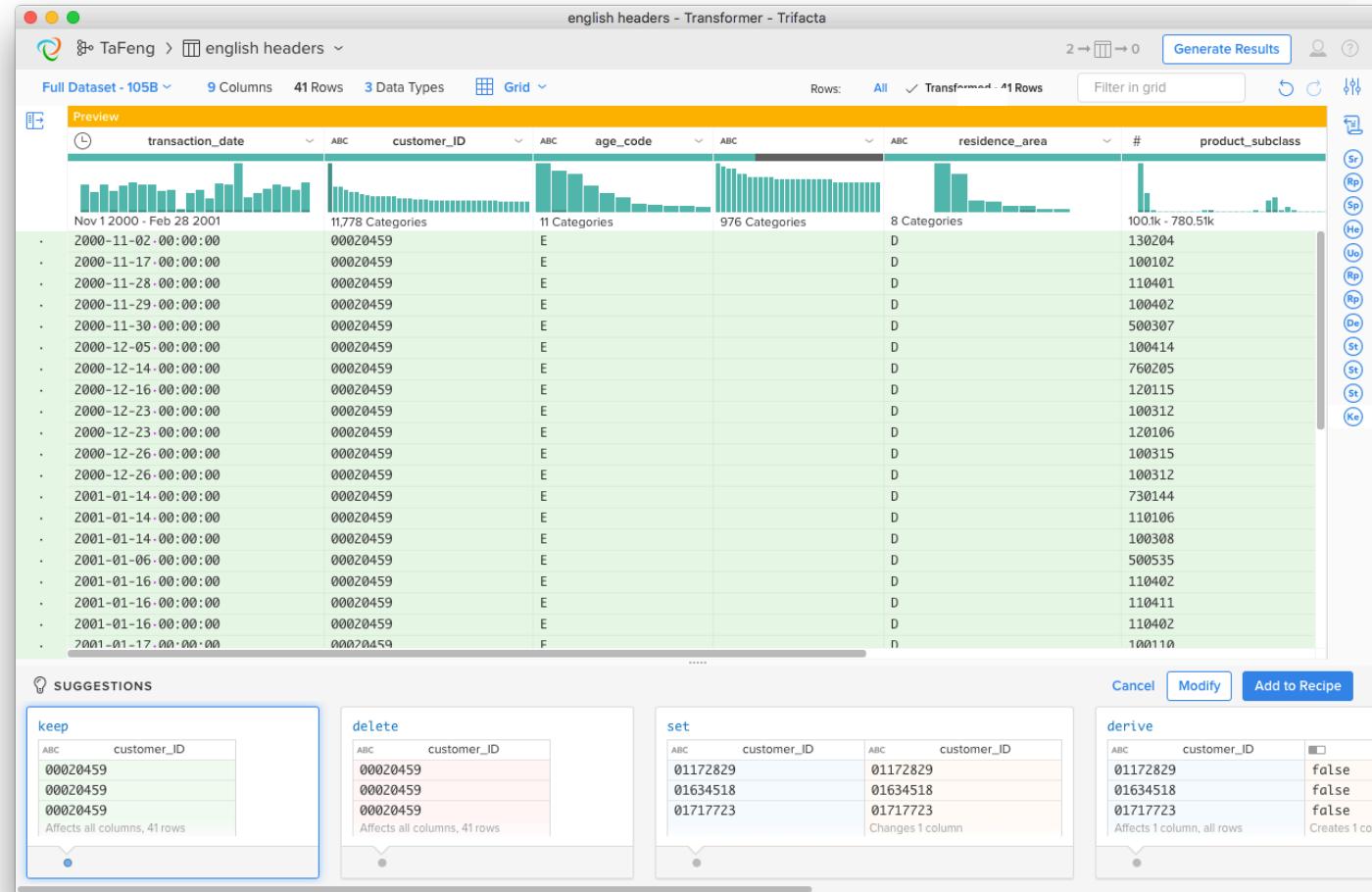
New Step   [Switch to editor](#)

Choose a transformation

[Choose transformation](#)

Cancel   [Add to Recipe](#)

# Assessing Faithfulness Within Records: Dependencies and Correlations



# 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: `customer_id`  $\rightarrow$  `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!

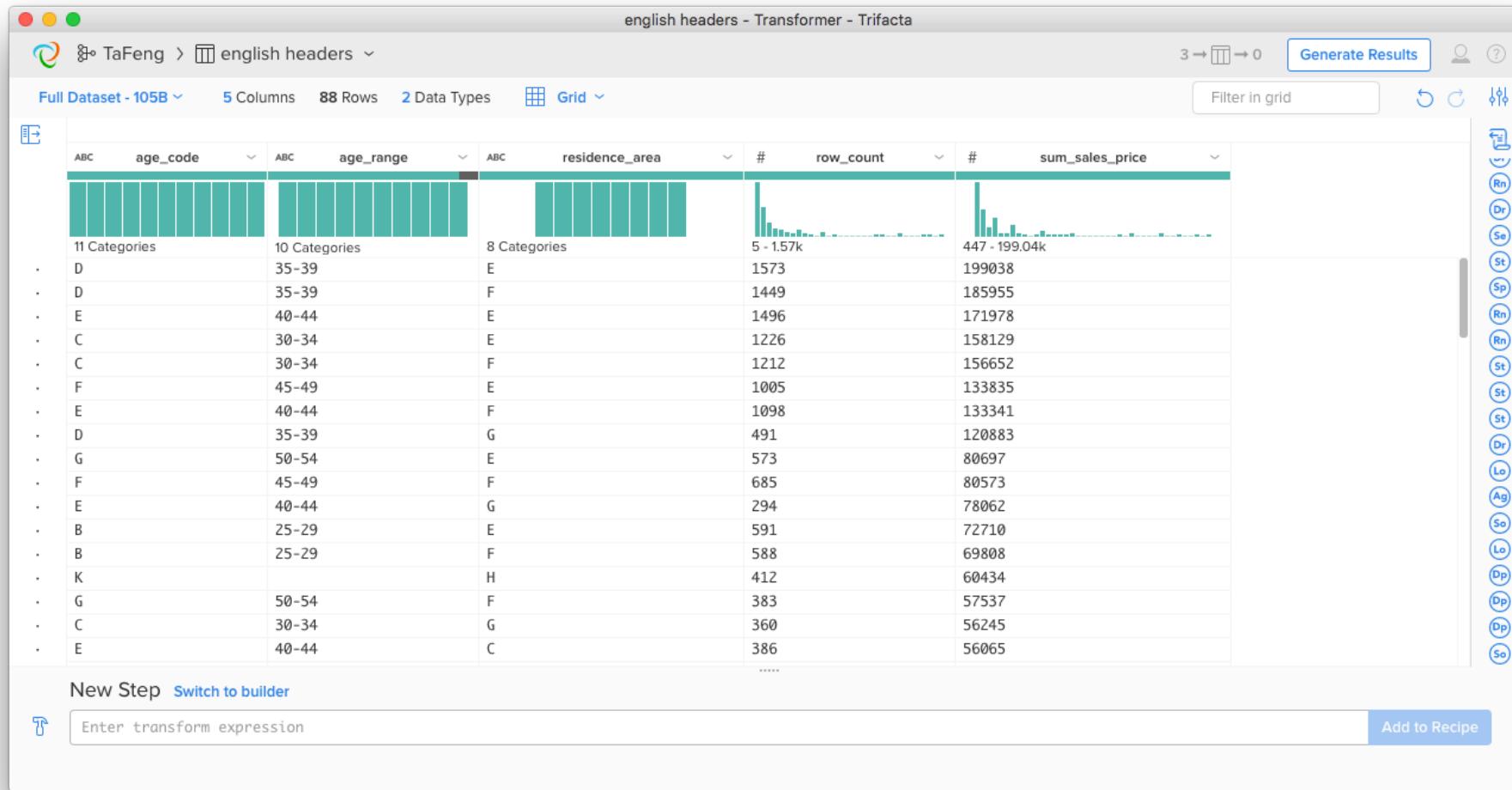
# Summing Up: Faithfulness

- 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



# What About Non-Rectangular Data

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

# Unnesting Nested Data Types

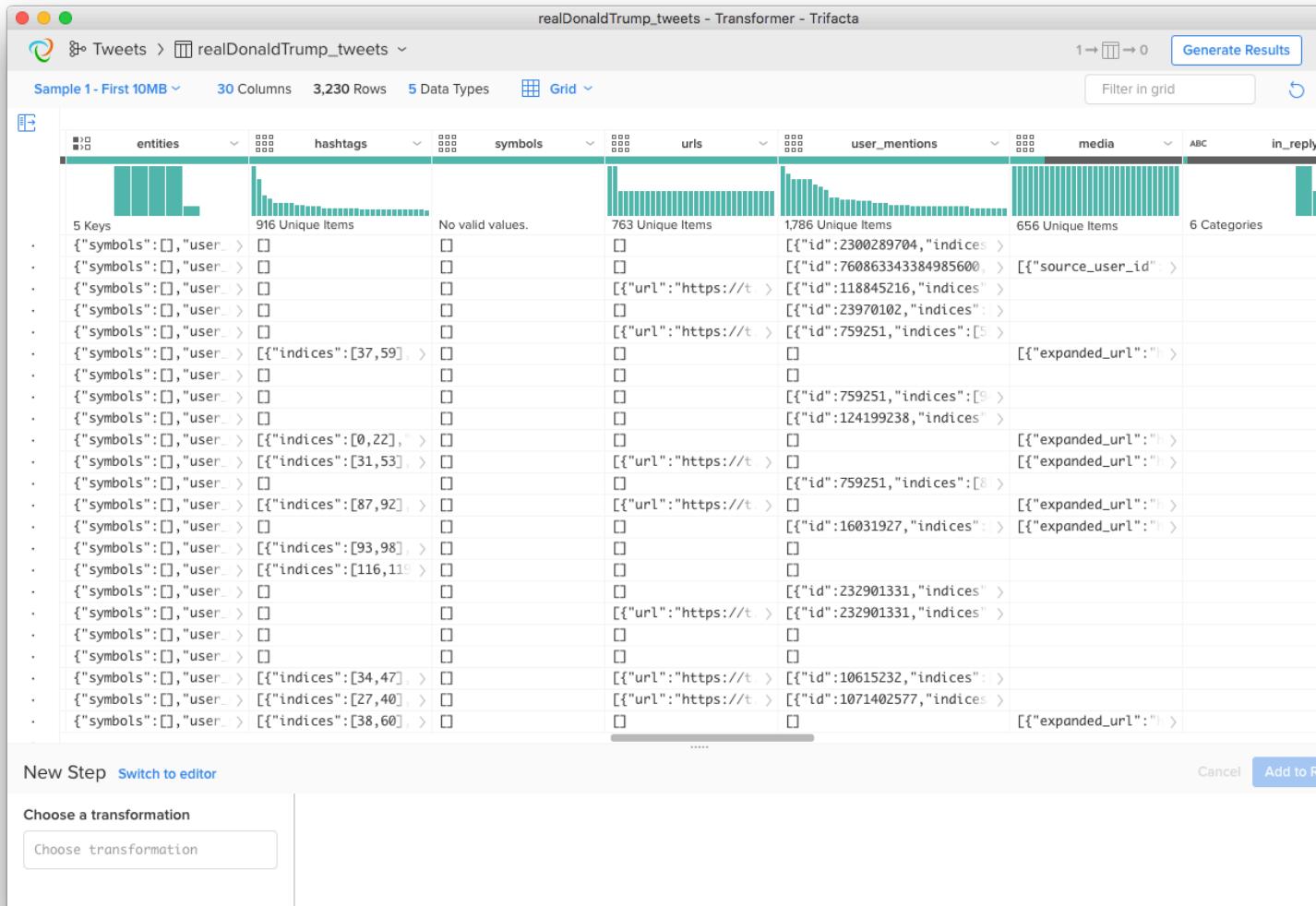
More on this in future lectures!

## Maps

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

## Arrays

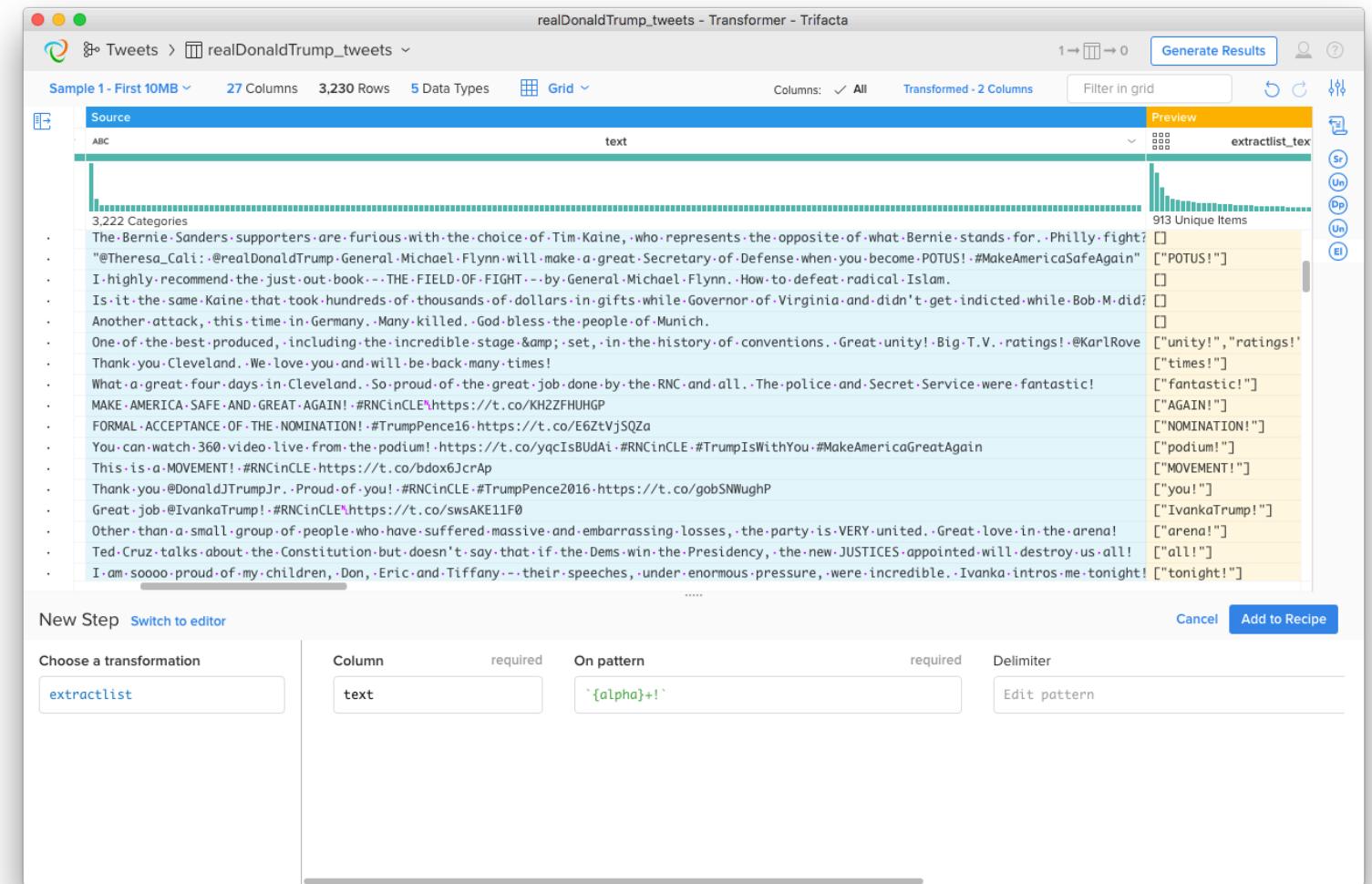
- Or lists



More on this in future  
lectures!

# Playing with Text I: String Manipulation

## Regular expression (regex) extraction



# 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		
Below is the list of postal codes in Taiwan.		
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
Keelung, New Taipei, Matsu, Yilan		
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
35
Affects all columns, 5056 rows

set
# age_min
35
30
40
# age_min
30
40
Changes 1 column

derive			
# age_min	column1		
35	true		
30	false		
40	false		
Affects 1 column, all rows		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!

# NSFW

Never Stop Fiercely Wrangling

Your analysis depends on it.

More fun:

<https://github.com/Quartz/bad-data-guide>

