#### Data Motivation

Sunday, July 31, 2016 1:56 PM

- Finding files with data
- Extracting raw data to be used
- Data typically does not look neat and organized
  - $\circ\ \ \mbox{ Raw file with data that must be parsed and analyzed to find meaning}$
- Goal:
  - o Raw data -> Processing script -> tidy data

#### Raw and Processed Data

Sunday, July 31, 2016 2:00 PM

- Raw data
  - o Original source of data
  - Hard for data analyses
  - Data analysis includes pr

Raw;

- Orig. source

- Hard for analysis

- May only need to be processed once

> processing is a step

Processedi

- Ready for analysis

Standards are available

Standards are available

recorded

rechnically - Stapes should be recorded

raw Processing Pipeline:

DNA Sequencing:

DNA Chunks -> colors attached

Burd on levely - take pictury

of shown color, and imany of

different product sequences product

# Components of Tidy Data

Sunday, July 31, 2016 2:06 PM

Raw Thy: Components:

- Raw dataset

- tily dataset

- Code book - each var t values in

meta Lata

- Explicit recording of process of

data

Kow dator 1. No saftwar on lota L. No manipulating 3. No removed data 4. No summerizing data Tidy Datas 1. Variable each column 2. Observations in rows 3. table per kind of variable 4. Multiple tables-> 17hk via column -tips: · Now at top w/ Var names · Variable names, human readable · One file per talk

Code book (Meta deta)

1. into about varslands!)

2. Summary choirs made

3. Study lesign used

Typically word/markdown

Instruction List:

- I deally a script,

no parametrus

input > van

output > tidy

If not I script, separate

into steps

# **Downloading Files**

Sunday, July 31, 2016 2:26 PM

- Best to do downloading via program language (eg R), so that step can be done via instruction script
- Get/set working directory
- Check for + create necessary directories
- Getting from Internet: download.file():
  - o Url
  - o Destfile (destination files)
  - Method (sometimes necessary to specify for secure connections)
- Record when downloaded for tracking purposes

# Reading Files (with R)

Sunday, July 31, 2016 2:33 PM

- Excel:
  - o Download/grab file
  - o Pretty straightforward from there... Google it
- XML:
  - Tags = general labels
    - Start/End tags: <section> </section>
    - Empty tags: <phrase />
  - o Elements are examples of tags, specific tags
    - <Greeting>
    - <Swag>
  - o Attributes are components of tags/labels
  - o XML Package
  - RootNode <- xmlRoot(doc)</li>
  - Names of root node:
    - xmlNames(RootNode)
    - Also can access like an array (rootNode[[1]])
    - xmlSApply()- apply across XML tree

#### Reading Databases

Sunday, July 31, 2016 3:56 PM

- MySQL
  - Open source database
  - o Package: RMySQL
  - o Method:
    - ucscDb <- dbConnect(MySQL(), user="",host="")</p>
    - Result <- dbGetQuery(ucscDb, "show databases;");</li>
    - Disconnect: dbDisconnect(ucscDb)
    - Can grab whole tables, or subsets via specific queries
    - Always close the connection when you don't need it anymore
- HDF5
  - Large datasets
  - Range of data types
  - Hierarchical data format
  - o Read in groups, different indices of tables
- From the Web
  - o Extract data from HTML code
  - Watch for breaching terms of service
  - o Reading too much quickly can get IP blocked
  - o Open connection
  - o Can use XML library, htmlparser to recognize general HTML tags
  - o GET(url)- Parse HTML, get info from the page
    - Websites with passwords
  - o Handles: root urls that can be built upon
- APIs
  - o Combine R with APIs for different apps to get data from these apps (eg Twitter)
  - After retrieving data paired to the software, process the data
  - o API documentation will show what interactions with server are available
  - Use httr package
- Other sources
  - There's an R package for almost anything
    - Databases: RPostressSQL, RODBC, Rmongo
    - Images: jpeg, readbitmap, png
    - GIS data
    - Music data: tuneR, seewave

# **Subsetting Tricks**

Sunday, July 31, 2016

11:24 PM

Missing valuessi

\* Luwich (\_\_) 

To returns correspondly

induces

Order X [order[---), ]

Orders rows
in terms of that

in terms of that

chind() = Bigling vows

### Summarizing Data

Sunday, July 31, 2016

11:32 PM

head= 6 top rows tall= bottom 3

Summary ():

- Text/factor vars-7 gin #

- Quantitative => give basic stat data

quantile();

( mín, max, max, max)

- Different dirts.

+ probs of quent data

t.bl({) "

-dutogen table for specified set

Check for missing value

sum (is. ha ())

or any (

all ( condition) =) don condition ever

Specitic Characterity

% in 1 =7

as data france ( )

(ross tabic

x tabs ( Var N Column to gration)

# Creating New Variables

Sunday, July 31, 2016 11:41 PM

Why?

Rew Lata dozsn't have needed value

Transforms necessary

Easy all to data frances

Common:

Missingness indicators

Lattley up quent
Transforms

• Sequences create indexed datasets

#### Dplyr

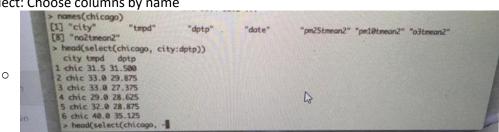
Monday, August 1, 2016 1:41 AM

dplyr = package for data frame doly rassumptionsi - 1 column = var/measure/cher. - Prinary imp. = R - Other laps good for relational databases

- data frame ) dp/yr-> new data maniputalism frame

#### Verbs:

Select: Choose columns by name



• Filter: subset rows based on conditions

```
> chic.f <- filter(chicago, pm25tmean2 > 30)
> head(chic.f, 10)
                                    date pm25tmean2 pm10tmean2 o3tmean2 no2tmean2
    city tmpd dptp
 1 chic 23 21.9 1998-01-17 38.10 32.46154 3.180556 25.30000
                                               33.95 38.69231 1.750000 29.37630
 2 chic 28 25.8 1998-01-23 33.95 38.69231 1.750000 29.37630
3 chic 55 51.3 1998-04-30 39.40 34.00000 10.786232 25.31310
 4 chic 59 53.7 1998-05-01 35.40 28.50000 14.295125 31.42905 5 chic 57 52.0 1998-05-02 33.30 35.00000 20.662879 26.79861 6 chic 57 56.0 1998-05-07 32.10 34.50000 24.270422 33.99167 7 chic 75 65.8 1998-05-15 56.50 91.00000 38.573007 29.03261
   7 chic
                                             33.80 26.00000 17.890810 25.49668
30.30 64.50000 37.018865 37.93056
41.40 75.00000 40.080902 32.59054
   8 chic 61 59.0 1998-06-09
   9 chic 73 60.3 1998-07-13
   10 chic 78 67.1 1998-07-14 41.40 75.00000 40.08 > chic.f <- filter(chicago, pm25tmean2 > 30 & tmpd > 80)
    > head(chic.f)
                                    date pm25tmean2 pm10tmean2 o3tmean2 no2tmean2
     city trpd dptp date
1 chic 81 71.2 1998-08-23
                                                                  59.0 45.86364 14.32639
                                                39.6000
                                                                   50.5 50.66250 20.31250
     2 chic 81 70.4 1998-09-06
                                                31.5000
                                                                  58.5 33.00380 33.67500
                82 72.2 2001-07-20
84 72.9 2001-08-01
                                                32.3000
     3 chic
4 chic
                                                                  81.5 45.17736 27.44239
70.0 37.98047 27.62743
                                                43.7000
```

- Arrange: reorder rows of data frame
- Rename: rename variable to something more intuitive
- Mutate: transform existing vars to make new ones
- Groupby: split into categorical data
- Summarize: self explanatory
- %>%: pipeline operator, feed into the right

#### **Regular Expressions**

Monday, August 1, 2016 4:26 AM

- Literals + metacharacters
- Literal = actual word
- ^(phrase) = phrase at beginning of line
- (phrase)\$ = phrase at end of line
- [Character class]- set of chars acceptable in match (eg [Bb] takes B or b)
- Range of characters:
  - [a-z] = all lowercase
  - [a-zA-Z]
  - o [0-9]
- ^ inside a char class indicates not what's in the char class
  - Eg [^?.]\$ = not ? Or . At end
- . = any character
- | = or
- ? = after any expression, that it's optional
- \ = precede other metacharacters when used as liberals
- \* = any number including none of item
- + = any number including 1
- {a, b} = indicates interval, i.e. The maximum and minimum occurrences
- \number = repetitions of found text expected