# Nuts and Bolts

Matthew Gentzkow
Jesse M. Shapiro

Chicago Booth and NBER

- We have focused on the statistical / econometric issues that arise with big data
- In the time that remains, we want to spend a little time on the *practical* issues...

- We have focused on the statistical / econometric issues that arise with big data
- In the time that remains, we want to spend a little time on the *practical* issues...
  - E.g., where do you actually put a 2 TB dataset?

- We have focused on the statistical / econometric issues that arise with big data
- In the time that remains, we want to spend a little time on the *practical* issues...
    - E.g., where do you actually put a 2 TB dataset?
- Goal: Sketch some basic computing ideas relevant to working with large datasets.

- We have focused on the statistical / econometric issues that arise with big data
- In the time that remains, we want to spend a little time on the *practical* issues...
    - E.g., where do you actually put a 2 TB dataset?
- Goal: Sketch some basic computing ideas relevant to working with large datasets.
- Caveat: We are all amateurs

# The Good News

- Much of what we've talked about here you can do on your laptop
  - Your OS knows how to do parallel computing (multiple processors, multiple cores)
  - Many "big" datasets are $< 5$ GB
  - Save the data to local disk, fire up Stata or R, and off you go...

# How Big is Big?

| Congressional record text (1870-2010) | ≈50 GB |
|---|---|
| Congressional record pdfs (1870-2010) | ≈500 GB |
| Nielsen scanner data (34k stores, 2004-2010) | ≈5 TB |
| Wikipedia (2013) | ≈6 TB |
| 20% Medicare claims data (1997-2009) | ≈10 TB |
| Facebook (2013) | ≈100,000 TB |
| All data in the world | ≈2.7 billion TB |

# Outline

- Software engineering for economists

- Databases
- Cluster computing
- Scenarios

# Software Engineering for Economists

# Motivation

- A lot of the time spent in empirical research is writing, reading, and debugging code.
- Common situations...

# Broken Code

# Incoherent Data

# Rampant Duplication

# Replication Impossible

# Tons of Versions

# This Talk

- We are not software engineers or computer scientists.
- But we have learned that most common problems in social sciences have analogues in these fields and there are standard solutions.
- Goal is to highlight a few of these that we think are especially valuable to researchers.
- Focus on incremental changes: one step away from common practice.

# Automation

# Raw Data

Data from original source...

Left spreadsheet (tab: chips), cell D2 = 1012:

| county | state | year | chip_sales |
|--------|-------|------|-----------|
| Autauga | AL | 1940 | 1012 |
| Autauga | AL | 1941 | 1020 |
| Autauga | AL | 1942 | 1034 |
| Autauga | AL | 1943 | 1058 |
| Autauga | AL | 1944 | 1085 |
| Autauga | AL | 1945 | 1148 |
| Autauga | AL | 1946 | 1205 |
| Autauga | AL | 1947 | 1287 |
| Autauga | AL | 1948 | 1299 |
| Autauga | AL | 1949 | 1344 |
| Autauga | AL | 1950 | 1365 |
| Autauga | AL | 1951 | 1397 |
| Autauga | AL | 1952 | 1455 |
| Autauga | AL | 1953 | 1501 |
| Autauga | AL | 1954 | 1582 |
| Autauga | AL | 1955 | 1656 |
| Autauga | AL | 1956 | 1723 |
| Autauga | AL | 1957 | 1795 |
| Autauga | AL | 1958 | 1878 |

Right spreadsheet (tab: tv), cell C32:

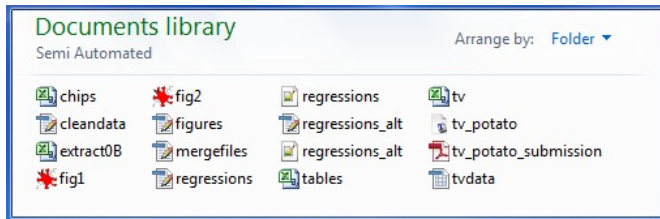| county | state | year_tv_introduced |
|--------|-------|--------------------|
| Autauga | AL | 1940 |
| Baldwin | AL | 1935 |
| Barbour | AL | 1942 |
| Bibb | AL | 1942 |
| Blount | AL | 1939 |
| Bullock | AL | 1945 |
| Butler | AL | 1942 |
| Calhoun | AL | 1936 |
| Chambers | AL | 1940 |
| Cherokee | AL | 1939 |
| Chilton | AL | 1941 |
| Choctaw | AL | 1942 |
| Clarke | AL | 1940 |
| Clay | AL | 1941 |
| Cleburne | AL | 1943 |
| Coffee | AL | 1936 |
| Colbert | AL | 1937 |
| Conecuh | AL | 1940 |
| Coosa | AL | 1943 |

# Manual Approach

- Open spreadsheet
- Output to text files
- Open Stata
- Load data, merge files
- Compute log(chip sales)
- Run regression
- Copy results to MS Word and save

# Manual Approach

- Two main problems with this approach
  - Replication: how can we be sure we'll find our way back to the exact same numbers?
  - Efficiency: what happens if we change our mind about the right specification?

# Semi-automated Approach



- Problems
  - Which file does what?
  - In what order?

# Fully Automated Approach

```
File:  rundirectory.bat
stattransfer export_to_csv.stc
statase -b mergefiles.do
statase -b cleandata.do
statase -b regressions.do
statase -b figures.do
pdflatex tv_potato.tex
```
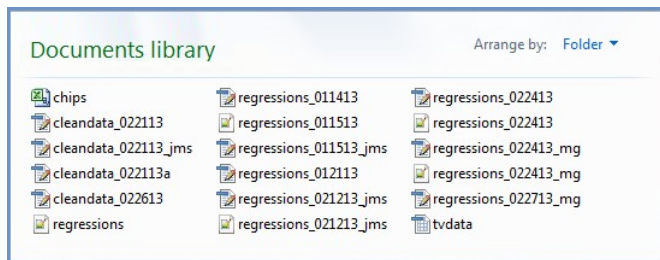
- All steps controlled by a shell script
- Order of steps unambiguous
- Easy to call commands from different packages

# Make

- Framework to go from source to target
- Tracks dependencies and revisions
- Avoids rebuilding components that are up to date
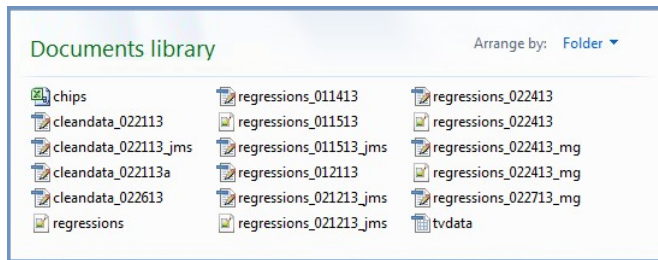- Used to build executable files

# Version Control

Documents library — Arrange by: Folder ▾

chips
cleandata_022113
cleandata_022113_jms
cleandata_022113a
cleandata_022613
regressions

regressions_011413
regressions_011513
regressions_011513_jms
regressions_012113
regressions_021213_jms
regressions_021213_jms

regressions_022413
regressions_022413
regressions_022413_mg
regressions_022413_mg
regressions_022713_mg
tvdata

- Dates demarcate versions, initials demarcate authors
- Why do this?
  - Facilitates comparison
  - Facilitates "undo"

- Why not do this?
  - It's a pain: always have to remember to "tag" every new file
  - It's confusing:
    - Which log file came from `regressions_022713_mg.do`?
    - Which version of `cleandata.do` makes the data used by `regressions_022413.do`?
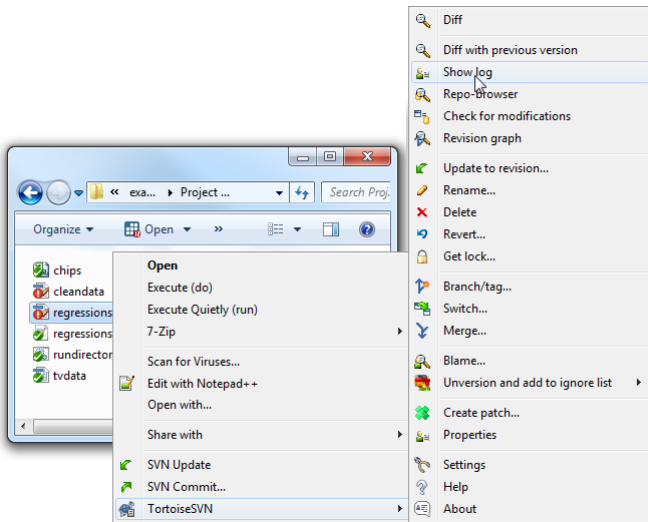  - It fails the market test: No software firm does it this way

# Version Control

- Software that sits "on top" of your filesystem
  - Keeps track of multiple versions of the same file
  - Records date, authorship
  - Manages conflicts

- Benefits
  - Single authoritative version of the directory
  - Edit without fear: an undo command for everything

Documents library

- chips
- cleandata
- regressions
- regressions
- rundirectory
- tvdata

# Life After Version Control

# Life After Version Control

# Life After Version Control

Documents library

- chips
- cleandata
- regressions
- regressions
- rundirectory
- tvdata

- Aside: If you always run rundirectory.bat before you commit, you guarantee replicability.

# Directories

- Pros: Self-contained, simple
- Cons:
  - Have to rerun everything for every change
  - Hard to figure out dependencies

# Functional Directories



Documents library

- build
  - code
    - cleandata
    - mergefiles
    - rundirectory
  - input
    - extract0B
  - output
    - tvdata
  - temp
    - chips
    - tv
- analysis
  - code
    - figures
    - regressions
    - regressions_alt
    - rundirectory
    - getinput
  - input
    - tvdata
  - output
    - fig1
    - fig2
    - tables
  - temp
    - regressions
    - regressions_alt

# Keys

# Research Assistant Output

| county | state | cnty_pop | state_pop | region |
|--------|-------|----------|-----------|--------|
| 36037 | NY | 3817735 | 43320903 | 1 |
| 36038 | NY | 422999 | 43320903 | 1 |
| 36039 | NY | 324920 | . | 1 |
| 36040 | . | 143432 | 43320903 | 1 |
| . | NY | . | 43320903 | 1 |
| 37001 | VA | 3228290 | 7173000 | 3 |
| 37002 | VA | 449499 | 7173000 | 3 |
| 37003 | VA | 383888 | 7173000 | 4 |
| 37004 | VA | 483829 | 7173000 | 3 |

# Causes for Concern

| county | state | cnty_pop | state_pop | region |
|--------|-------|----------|-----------|--------|
| 36037  | NY    | 3817735  | 43320903  | 1      |
| 36038  | NY    | 422999   | 43320903  | 1      |
| 36039  | NY    | 324920   | .         | 1      |
| 36040  | .     | 143432   | 43320903  | 1      |
| .      | NY    | .        | 43320903  | 1      |
| 37001  | VA    | 3228290  | 7173000   | 3      |
| 37002  | VA    | 449499   | 7173000   | 3      |
| 37003  | VA    | 383888   | 7173000   | 4      |
| 37004  | VA    | 483829   | 7173000   | 3      |

# Relational Databases

| county | state | population |
|--------|-------|-----------|
| 36037 | NY | 3817735 |
| 36038 | NY | 422999 |
| 36039 | NY | 324920 |
| 36040 | NY | 143432 |
| 37001 | VA | 3228290 |
| 37002 | VA | 449499 |
| 37003 | VA | 383888 |
| 37004 | VA | 483829 |

| state | population | region |
|-------|-----------|--------|
| NY | 43320903 | 1 |
| VA | 7173000 | 3 |

- Each *variable* is an attribute of an *element* of the table
- Each table has a *key*
- Tables are connected by *foreign keys* (state field in the county table)
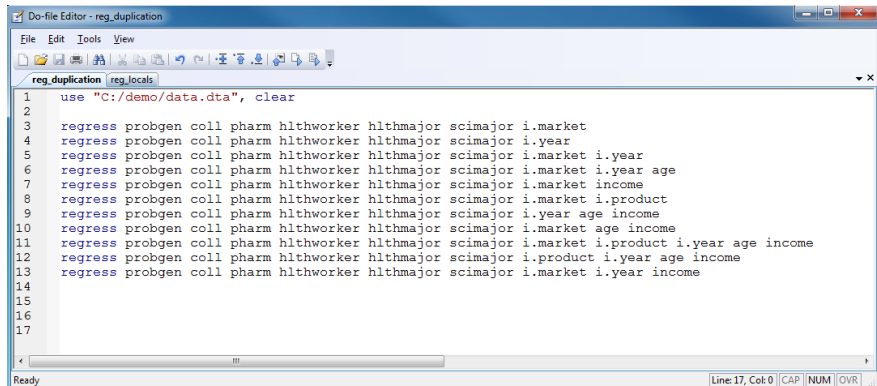
# Steps

- Store data in normalized format as above
  - Can use flat files, doesn't have to be fancy relational database software
- Construct a second set of files with key transformations
  - e.g., log population
- Merge data together and run analysis

- What to do with enormous databases?

# Abstraction

# Rampant Duplication
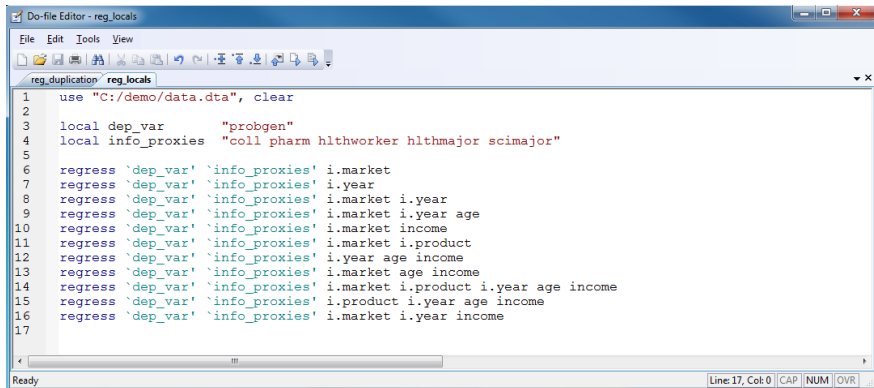


```
use "C:/demo/data.dta", clear

regress probgen coll pharm hlthworker hlthmajor scimajor i.market
regress probgen coll pharm hlthworker hlthmajor scimajor i.year
regress probgen coll pharm hlthworker hlthmajor scimajor i.market i.year
regress probgen coll pharm hlthworker hlthmajor scimajor i.market i.year age
regress probgen coll pharm hlthworker hlthmajor scimajor i.market income
regress probgen coll pharm hlthworker hlthmajor scimajor i.market i.product
regress probgen coll pharm hlthworker hlthmajor scimajor i.market age income
regress probgen coll pharm hlthworker hlthmajor scimajor i.market i.product i.year age income
regress probgen coll pharm hlthworker hlthmajor scimajor i.year age income
regress probgen coll pharm hlthworker hlthmajor scimajor i.product i.year age income
regress probgen coll pharm hlthworker hlthmajor scimajor i.market i.year income
```

# Abstracted



```
1    use "C:/demo/data.dta", clear
2
3    local dep_var      "probgen"
4    local info_proxies "coll pharm hlthworker hlthmajor scimajor"
5
6    regress `dep_var' `info_proxies' i.market
7    regress `dep_var' `info_proxies' i.year
8    regress `dep_var' `info_proxies' i.market i.year
9    regress `dep_var' `info_proxies' i.market i.year age
10   regress `dep_var' `info_proxies' i.market income
11   regress `dep_var' `info_proxies' i.market i.product
12   regress `dep_var' `info_proxies' i.year age income
13   regress `dep_var' `info_proxies' i.market age income
14   regress `dep_var' `info_proxies' i.market i.product i.year age income
15   regress `dep_var' `info_proxies' i.product i.year age income
16   regress `dep_var' `info_proxies' i.market i.year income
17
```

# Three Leave-Out Means



```stata
1    * Per capita consumption within state
2    egen total_pc_potato = total(pc_potato), by(state)
3    egen total_obs = count(pc_potato), by(state)
4    gen leaveout_state_pc_potato = (total_pc_potato - pc_potato)/(total_obs - 1)
5
6    * Per capita consumption within metro area
7    egen total_pc_potato = total(pc_potato), by(metroarea)
8    egen total_obs = count(pc_potato), by(state)
9    gen leaveout_metro_pc_potato = (total_pc_potato - pc_potato)/(total_obs - 1)
10
11   * Per household consumption within metro area
12   egen total_hh_potato = total(hh_potato), by(metroarea)
13   egen total_obs = count(hh_potato), by(state)
14   gen leaveout_metro_hh_potato = (total_hh_potato - pc_potato)
15
```

# Copy and Paste Errors



```
 1    * Per capita consumption within state
 2    egen total_pc_potato = total(pc_potato), by(state)
 3    egen total_obs = count(pc_potato), by(state)
 4    gen leaveout_state_pc_potato = (total_pc_potato - pc_potato)/(total_obs - 1)
 5
 6    * Per capita consumption within metro area
 7    egen total_pc_potato = total(pc_potato), by(metroarea)
 8    egen total_obs = count(pc_potato), by(state )
 9    gen leaveout_metro_pc_potato = (total_pc_potato - pc_potato)/(total_obs - 1)
10
11    * Per household consumption within metro area
12    egen total_hh_potato = total(hh_potato), by(metroarea)
13    egen total_obs = count(hh_potato), by(state )
14    gen leaveout_metro_hh_potato = (total_hh_potato - pc_potato )
15
```

# Abstracted

```
 1
 2   program leaveout_mean
 3       syntax, invar(varname) outvar(name) byvar(varname)
 4       tempvar tot_invar count_invar
 5       egen `tot_invar' = total(`invar'), by(`byvar')
 6       egen `count_invar' = count(`invar'), by(`byvar')
 7       gen `outvar' = (`tot_invar' - `invar') / (`count_invar' - 1)
 8   end
 9
10   leaveout_mean, invar(pc_potato) outvar(leaveout_state_pc_potato) byvar(state)
11   leaveout_mean, invar(pc_potato) outvar(leaveout_metro_pc_potato) byvar(metro)
12   leaveout_mean, invar(hh_potato) outvar(leaveout_metro_hh_potato) byvar(metro)
13
```

# Documentation

# Too Much Documentation

# Too Much Documentation

# Too Much Documentation

# Too Much Documentation



```
/***********************************************************
run_regressions.do fits our county fixed effects model with and
without controls for ranch dip sales and salsa consumption
***********************************************************/

use "C:\demo\data\chips_tv_1940.dta", clear

forval i=1941/2012 {
    append using "C:\demo\data\chips_tv_`i'.dta"
}

save "C:\demo\data\chips_tv_allyears.dta", replace
```

# Too Much Documentation

# Unclear Code



```
1    local el=0.4/0.2
2    compwlf, input(`el')
3
4
5
```

# Self-Documenting Code



Do-file Editor - unclear_code

File  Edit  Tools  View

unclear_code

```
1    local el=0.4/0.2
2    compwlf, input(`el')
3
4
5
```

Ready                                                    Line: 5, Col: 0  CAP  NUM  OVR



Do-file Editor - selfdoc_code

File  Edit  Tools  View

selfdoc_code

```
1    local percent_change_in_quantity = -0.4
2    local percent_change_in_price = 0.2
3    local elasticity = `percent_change_in_quantity'/`percent_change_in_price'
4    compute_welfare_loss, elasticity(`elasticity')
5
```

Ready                                                    Line: 5, Col: 0  CAP  NUM  OVR

# Management

# A Friendly Chat



Hey Matt,

Do you have that robustness check where we control for the amount of ranch dip sold in each county? I am writing the section on dipping sauces and wanted to mention it.

Jesse

# A Friendly Chat



potato chips - Message (HTML)

File    Message    Insert    Options    Format Text    Review

To...    Shapiro, Jesse

Cc...

Subject:    potato chips

Sorry, I thought you were doing that because it's similar to that other thing you were doing with controlling for salsa sales. Let me know if you want to do it or if you want me to take over.

MG

See more about: Shapiro, Jesse.

# A Friendly Chat



potato chips - Message (HTML)

File    Message    Insert    Options    Format Text    Review

Calibri (Body)  16

To...     Gentzkow, Matthew
Cc...     Sinkinson, Michael
Subject:  potato chips

I thought Matt was doing ranch dip and Mike was doing salsa?

Jesse

See more about: Gentzkow, Matthew.

# A Friendly Chat

potato chips - Message (HTML)

File | Message | Insert | Options | Format Text | Review

Paste | Cut | Copy | Format Painter
Clipboard

Calibri (Body) | 11 | **B** *I* <u>U</u>
Basic Text

Address Book | Check Names
Names

Attach File | Attach Item | Signature
Include

Follow Up
High Importance
Low Importance
Tags

Zoom
Zoom

To... | Shapiro, Jesse; Gentzkow, Matthew

Cc...

Subject: | potato chips

I did the salsa robustness check two weeks ago. See my e-mail from 8/14, 9:36am.

Mike

See more about: Shapiro, Jesse.

# A Friendly Chat



To... Sinkinson, Michael; Gentzkow, Matthew

Subject: potato chips

Right, but in that e-mail you were controlling for the log of salsa consumption. I thought we agreed we wanted the level of consumption?

Jesse

See more about: Sinkinson, Michael.

# A Friendly Chat

# Task Management

# Parting Thoughts

# Code and Data

- Data are getting larger
- Research is getting more collaborative
- Need to manage code and data responsibly for collaboration and replicability
- Learn from the pros, not from us

# Databases

# What is a Database?

- Database Theory
  - Principles for how to store / organize / retrieve data efficiently (normalization, indexing, optimization, etc.)
- Database Software
  - Manages storage / organization / retrieval of data (SQL, Oracle, Access, etc.)
  - Economists rarely use this software because we typically store data in flat files & interact with them using statistical programs
  - When we receive extracts from large datasets (the census, Medicare claims, etc.) someone else often interacts with the database on the back end

# Normalization

- *"Database Normalization* is the process of organizing the fields and tables of a relational database to minimize redundancy and dependency. Normalization usually involves dividing large tables into smaller (and less redundant) tables and defining relationships between them."

# Benefits of Normalization

- Efficient storage
- Efficient modification
- **Guarantees coherence**
- **Makes logical structure of data clear**

- Medicare claims data for 1997-2010 are roughly 10 TB
- These data are stored at NBER in thousands of zipped SAS files

# Indexing

- Medicare claims data for 1997-2010 are roughly 10 TB
- These data are stored at NBER in thousands of zipped SAS files
- To extract, say, all claims for heart disease patients aged 55-65, you would need to read every line of every one of those files
  - THIS IS SLOW!!!

- The obvious solution, long understood for book, libraries, economics journals, and so forth, is to build an index
- Database software handles this automatically
  - Allows you to specify fields that will be often used for lookups, subsetting, etc. to be indexed
  - For the Medicare data, we could index age, gender, type of treatment, etc. to allow much faster extraction

- Benefits
  - Fast lookups
  - Easy to police data constraints

- Costs
  - Storage
  - Time

- Database *optimization* is the art of tuning database structure and indexing for a specific set of needs

- Traditional databases are optimized for *operational* environments
  - Bank transactions
  - Airline reservations
  - etc.

# Data Warehouses

- Traditional databases are optimized for *operational* environments
  - Bank transactions
  - Airline reservations
  - etc.

- Characteristics
  - Many small reads and writes
  - Many users accessing simultaneously
  - Premium on low latency
  - Only care about current state

- In analytic / research environments, however, the requirements are different
  - Frequent large reads, infrequent writes
  - Relatively little simultaneous access
  - Value throughput relative to latency
  - May care about history as well as current state
  - Need to create and re-use many custom extracts

# Data Warehouses

- In analytic / research environments, however, the requirements are different
  - Frequent large reads, infrequent writes
  - Relatively little simultaneous access
  - Value throughput relative to latency
  - May care about history as well as current state
  - Need to create and re-use many custom extracts
- Database systems tuned to these requirements are commonly called "data warehouses"

# Distributed Computing

- Definition: Computation shared among many independent processors

# Distributed Computing

- Definition: Computation shared among many independent processors
- Terminology
  - Distributed vs. Parallel (latter usually refers to systems with shared memory)
  - Cluster vs. Grid (latter usually more decentralized & heterogeneous)

# On Your Local Machine

- Your OS can run multiple processors each with multiple cores
- Your video card has hundreds of cores
- Stata, R, Matlab, etc. can all exploit these resources to do parallel computing

# On Your Local Machine

- Your OS can run multiple processors each with multiple cores
- Your video card has hundreds of cores
- Stata, R, Matlab, etc. can all exploit these resources to do parallel computing
- Stata
  - Buy appropriate "MP" version of Stata
  - Software does the rest

# On Your Local Machine

- Your OS can run multiple processors each with multiple cores
- Your video card has hundreds of cores
- Stata, R, Matlab, etc. can all exploit these resources to do parallel computing
- Stata
  - Buy appropriate "MP" version of Stata
  - Software does the rest
- R / Matlab
  - Install appropriate add-ins (*parallel* package in R, "parallel computing toolbox" in Matlab)
  - Include parallel commands in code (e.g., *parfor* in place of *for* in Matlab)

- Resources abound
  - University / department computing clusters
  - Non-commercial scientific computing grids (e.g., XSEDE)
  - Commercial grids (e.g., Amazon EC2)

- Resources abound
  - University / department computing clusters
  - Non-commercial scientific computing grids (e.g., XSEDE)
  - Commercial grids (e.g., Amazon EC2)
- Most of these run Linux w/ distribution handled by a "batch scheduler"
- Write code using your favorite application, then send it to scheduler with a bash script

- MapReduce is a programming model that facilitates distributed computing
  - Developed by Google around 2004, though ideas predate that

# MapReduce

- MapReduce is a programming model that facilitates distributed computing
  - Developed by Google around 2004, though ideas predate that
- Most algorithms for distributed data processing can be represented in two steps
  - **Map**: Process individual "chunk" of data to generate an intermediate "summary"
  - **Reduce:** Combine "summaries" from different chunks to produce a single output file

- MapReduce is a programming model that facilitates distributed computing
  - Developed by Google around 2004, though ideas predate that

- Most algorithms for distributed data processing can be represented in two steps
  - **Map**: Process individual "chunk" of data to generate an intermediate "summary"
  - **Reduce:** Combine "summaries" from different chunks to produce a single output file

- If you structure your code this way, MapReduce software will handle all the details of distribution:
  - Partitioning data
  - Scheduling execution across nodes
  - Managing communication between machines
  - Handling errors / machine failures

- Count words in a large collection of documents
  - Map: Document $i \rightarrow$ Set of (*word*, *count*) pairs $C_i$
  - Reduce: Collapse $\{C_i\}$, summing *count* within *word*

- Count words in a large collection of documents
  - Map: Document $i \rightarrow$ Set of ($word$, $count$) pairs $C_i$
  - Reduce: Collapse $\{C_i\}$, summing $count$ within $word$
- Extract medical claims for 65-year old males
  - Map: Record set $i \rightarrow$ Subset of $i$ that are 65-year old males $H_i$
  - Reduce: Append elements of $\{H_i\}$

- Count words in a large collection of documents
  - Map: Document $i \to$ Set of ($word, count$) pairs $C_i$
  - Reduce: Collapse $\{C_i\}$, summing $count$ within $word$

- Extract medical claims for 65-year old males
  - Map: Record set $i \to$ Subset of $i$ that are 65-year old males $H_i$
  - Reduce: Append elements of $\{H_i\}$

- Compute marginal regression for text analysis (e.g., Gentzkow & Shapiro 2010)
  - Map: Counts $x_{ij}$ of phrase $j \to$ Parameters $\left(\hat{\alpha}_j, \hat{\beta}_j\right)$ from $E\left(x_{ij}|y_i\right) = \alpha_j + \beta_j x_{ij}$
  - Reduce: Append $\left\{\hat{\alpha}_j, \hat{\beta}_j\right\}$

# MapReduce: Implementation

- MapReduce is the original software developed by Google
- Hadoop is the open-source version most people use (developed by Apache)
- Amazon has a hosted implementation (Amazon EMR)

# MapReduce: Implementation

- MapReduce is the original software developed by Google
- Hadoop is the open-source version most people use (developed by Apache)
- Amazon has a hosted implementation (Amazon EMR)
- How does it work?
  - Write your code as two functions called *map* and *reduce*
  - Send code & data to scheduler using bash script

# Distributed File Systems

- Data transfer is the main bottleneck in distributed systems
- For big data, it makes sense to distribute data as well as computation
  - Data broken up into chunks, each of which lives on a separate node
  - File system keeps track of where the pieces are and allocates jobs so computation happens "close" to data whenever possible

# Distributed File Systems

- Data transfer is the main bottleneck in distributed systems
- For big data, it makes sense to distribute data as well as computation
  - Data broken up into chunks, each of which lives on a separate node
  - File system keeps track of where the pieces are and allocates jobs so computation happens "close" to data whenever possible
- Tight coupling between MapReduce software and associated file systems
  - MapReduce $\rightarrow$ Google File System (GFS)
  - Hadoop $\rightarrow$ Hadoop Distributed File System (HDFS)
  - Amazon EMR $\rightarrow$ Amazon S3

# Distributed File Systems

# Scenarios

- *My data is* 100 *gb or less*

- *My data is* 100 *gb or less*
- Advice
  - Store data locally in flat files (csv, Stata, R, etc.)
  - Organize data in normalized tables for robustness and clarity
  - Run code serially or (if computation is slow) in parallel

# Scenario 2: Big Data, Small Analysis

- *My raw data is $> 100$ gb, but the extracts I actually use for analysis are $<< 100$ gb*

# Scenario 2: Big Data, Small Analysis

- *My raw data is $> 100$ gb, but the extracts I actually use for analysis are $<< 100$ gb*

- Example
  - Medicare claims data $\rightarrow$ analyze heart attack spending by patient by year
  - Nielsen scanner data $\rightarrow$ analyze average price by store by month

# Scenario 2: Big Data, Small Analysis

- *My raw data is $> 100$ gb, but the extracts I actually use for analysis are $<< 100$ gb*
- Example
  - Medicare claims data $\rightarrow$ analyze heart attack spending by patient by year
  - Nielsen scanner data $\rightarrow$ analyze average price by store by month
- Advice
  - Store data in relational database optimized to produce analysis extracts efficiently
  - Store extracts locally in flat files (csv, Stata, R, etc.)
  - Organize extracts in normalized tables for robustness and clarity
  - Run code serially or (if computation is slow) in parallel

# Scenario 2: Big Data, Small Analysis

- *My raw data is $> 100$ gb, but the extracts I actually use for analysis are $<< 100$ gb*
- Example
  - Medicare claims data $\rightarrow$ analyze heart attack spending by patient by year
  - Nielsen scanner data $\rightarrow$ analyze average price by store by month
- Advice
  - Store data in relational database optimized to produce analysis extracts efficiently
  - Store extracts locally in flat files (csv, Stata, R, etc.)
  - Organize extracts in normalized tables for robustness and clarity
  - Run code serially or (if computation is slow) in parallel
- Note: Gains to database increase for more structured data. For completely unstructured data, you may be better off using distributed file system + map reduce to create extracts.

- *My data is $> 100$ GB and my analysis code needs to touch all of the data*

- *My data is $> 100$ GB and my analysis code needs to touch all of the data*
- Example
  - 2 TB of SEC filing text $\rightarrow$ run variable selection using all data

- *My data is $> 100$ GB and my analysis code needs to touch all of the data*
- Example
  - 2 TB of SEC filing text $\rightarrow$ run variable selection using all data
- Advice
  - Store data in distributed file system
  - Use MapReduce or other distributed algorithms for analysis