

SOSC 4300/5500: Overview

Han Zhang

Feb 8, 2022

Outline

Computational Social Science

Logistics

Git and GitHub

Big Data: data acquisition

Big Data: Opportunities and Challenges

Before Digital Revolution

- And then we calculate some statistics from census surveys
- 1890 US census took **8 years** to clean and process by humans

DISTRICTS	Free white Males of 15 years and up, wards, including heads of families	Free white Males under fifteen years	Free white Females of 15 years and up, wards, including heads of families	All other free per- sons	Slaves	Total
Vermont	22435	22328	40505	255	16	85539
N. Hampshire	36080	34851	70160	630	158	141885
Maine	24384	24748	46870	538	NONE	96540
Massachusetts	95453	87289	190582	5401	NONE	378787
Rhode Island	16019	15799	32052	3407	948	68825
Connecticut	60523	54403	117448	2808	2761	137946
New York	83700	78122	151320	4654	21384	340180
New Jersey	45251	41416	83287	2762	11423	184139
Pennsylvania	110788	106948	206363	6537	3737	434373
Delaware	11783	12143	22384	3899	8887	59094
Maryland	55915	51339	101395	8043	103036	119728
Virginia	110936	116135	215046	12866	292627	747610
Kentucky	15154	17057	28922	114	12430	73677
N. Carolina	69988	77506	140710	4975	100571	393751
S. Carolina	35576	37722	60886	1801	107094	249973
Georgia	13103	14044	25739	398	29264	82548
	807094	791850	1541263	59150	694280	3893635
Total number of inhabitants of the United States exclusive of S. W. territory.						
	Free white Males of 21 years and upwards	Free white Males under 21 years of age	Free white Females	All other free persons	Slaves	Total
S. W. territory	6271	10277	15365	361	3417	35691
N. Ditto	—	—	—	—	—	—

With modern computers

- Invented in 1940s, modern personal computers become popular since 1980s

With modern computers

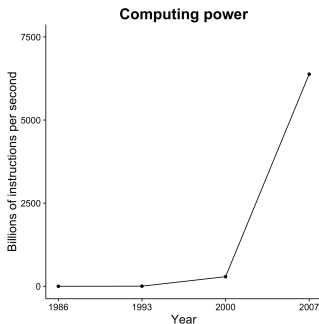
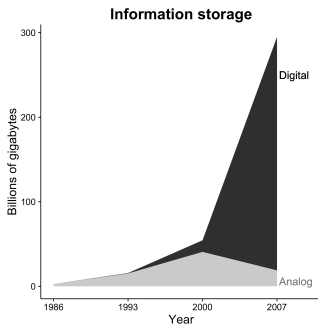
- Invented in 1940s, modern personal computers become popular since 1980s
- Calculation becomes much faster with computers

With modern computers

- Invented in 1940s, modern personal computers become popular since 1980s
- Calculation becomes much faster with computers
 - Imagine solving a regression by hand without computers!
- But data are still in analog format; they are represented in a physical way.

Welcome to the digital age

- Since 2000, both computing power and digital data are quickly increasing



Welcome to the digital age

- Computers everywhere, **digital traces** everywhere

Welcome to the digital age

- Computers everywhere, **digital traces** everywhere
 - personal computers, mobile phones, cars, watches, thermostats, CCTV cameras. . .

Welcome to the digital age

- Computers everywhere, **digital traces** everywhere
 - personal computers, mobile phones, cars, watches, thermostats, CCTV cameras. . .
 - these devices not only **calculate**; they also **measure** and **store** lots of digital data

Welcome to the digital age

- Computers everywhere, **digital traces** everywhere
 - personal computers, mobile phones, cars, watches, thermostats, CCTV cameras. . .
 - these devices not only **calculate**; they also **measure** and **store** lots of digital data
 - E.g., 20 years ago, you may walk into a bookstore and browse books; no traces will be left once you walk outside the book store
 - Now, your entire online browsing and purchasing behaviors are stored, and will be used for advertising or recommendation for similar products
 - Digital traces do not need to be on Internet!

Welcome to the digital age

- Computers everywhere, **digital traces** everywhere
 - personal computers, mobile phones, cars, watches, thermostats, CCTV cameras. . .
 - these devices not only **calculate**; they also **measure** and **store** lots of digital data
 - E.g., 20 years ago, you may walk into a bookstore and browse books; no traces will be left once you walk outside the book store
 - Now, your entire online browsing and purchasing behaviors are stored, and will be used for advertising or recommendation for similar products
 - Digital traces do not need to be on Internet!
 - E.g., octopus card swipes allow companies to locate your moving trajectories

Welcome to the digital age

- Computers everywhere, **digital traces** everywhere
 - personal computers, mobile phones, cars, watches, thermostats, CCTV cameras. . .
 - these devices not only **calculate**; they also **measure** and **store** lots of digital data
 - E.g., 20 years ago, you may walk into a bookstore and browse books; no traces will be left once you walk outside the book store
 - Now, your entire online browsing and purchasing behaviors are stored, and will be used for advertising or recommendation for similar products
 - Digital traces do not need to be on Internet!
 - E.g., octopus card swipes allow companies to locate your moving trajectories
- **Digital traces** are byproducts of peoples everyday actions, often collected by companies.

Welcome to the digital age

- Computers everywhere, **digital traces** everywhere
 - personal computers, mobile phones, cars, watches, thermostats, CCTV cameras. . .
 - these devices not only **calculate**; they also **measure** and **store** lots of digital data
 - E.g., 20 years ago, you may walk into a bookstore and browse books; no traces will be left once you walk outside the book store
 - Now, your entire online browsing and purchasing behaviors are stored, and will be used for advertising or recommendation for similar products
 - Digital traces do not need to be on Internet!
 - E.g., octopus card swipes allow companies to locate your moving trajectories
- **Digital traces** are byproducts of peoples everyday actions, often collected by companies.
 - Before digital age, they just fade away, but now they are kept

Big Data

- Together, we call digital traces and traditional data that are turned into digital data as **Big Data**
- “Big data are created and collected by companies and governments **for purposes other than research**”

Big data vs traditional social science data

- Traditional social science data are made for research
 - Although the data are small, they are ready to use for examining social science theories

Big data vs traditional social science data

- Traditional social science data are made for research
 - Although the data are small, they are ready to use for examining social science theories
- Big data are **repurposed** for research
 - They are big
 - But you need some effort to get what you want

Big data need different methods

- Previously social scientists have survey and sometimes small administrative data

Big data need different methods

- Previously social scientists have survey and sometimes small administrative data
 - Using various **regression** models to analyze the data

Big data need different methods

- Previously social scientists have survey and sometimes small administrative data
 - Using various **regression** models to analyze the data
- Big data are not only big, but also qualitatively different in their formats:

Big data need different methods

- Previously social scientists have survey and sometimes small administrative data
 - Using various **regression** models to analyze the data
- Big data are not only big, but also qualitatively different in their formats:
 - Texts

Big data need different methods

- Previously social scientists have survey and sometimes small administrative data
 - Using various **regression** models to analyze the data
- Big data are not only big, but also qualitatively different in their formats:
 - Texts
 - Images, Video, Audio
 - Networks

Big data need different methods

- Previously social scientists have survey and sometimes small administrative data
 - Using various **regression** models to analyze the data
- Big data are not only big, but also qualitatively different in their formats:
 - Texts
 - Images, Video, Audio
 - Networks
- Analyzing the above need other methods, in particular **machine learning**

Big data need different methods

- Previously social scientists have survey and sometimes small administrative data
 - Using various **regression** models to analyze the data
- Big data are not only big, but also qualitatively different in their formats:
 - Texts
 - Images, Video, Audio
 - Networks
- Analyzing the above need other methods, in particular **machine learning**
 - Regression in most times won't work!

Social scientists and data scientists

- Status quo:

Social scientists and data scientists

- Status quo:
- Social scientists: computational **social science**

Social scientists and data scientists

- Status quo:
- Social scientists: computational **social science**
 - Try to turn big data into small data, and then apply traditional regression models

Social scientists and data scientists

- Status quo:
- Social scientists: computational **social science**
 - Try to turn big data into small data, and then apply traditional regression models
- Data scientists: **computational** social science

Social scientists and data scientists

- Status quo:
- Social scientists: computational **social science**
 - Try to turn big data into small data, and then apply traditional regression models
- Data scientists: **computational** social science
 - Get more data, and apply some fancy machine learning algorithms over social data

Computational Social Science

- Social science itself is not enough, because data can only gets bigger

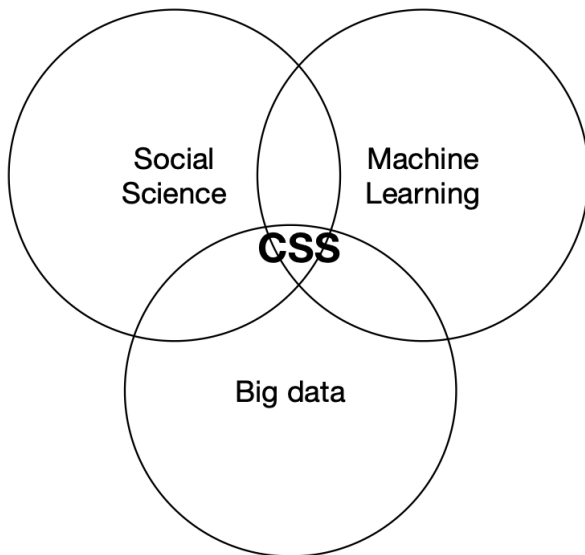
Computational Social Science

- Social science itself is not enough, because data can only gets bigger
- Data science itself is not enough, if we want to study social behaviors rigorously

Computational Social Science

- Social science itself is not enough, because data can only gets bigger
- Data science itself is not enough, if we want to study social behaviors rigorously
- Computational social science (CSS): bridging social and data science

Computational Social Science (CSS)



Study Goals

1. Describe the opportunities and challenges of computational social science
2. Evaluate computational social science research on social phenomena
3. Practice the essential techniques to analyze social big data, especially text data (covered in Tutorials)
 - Getting data
4. Propose research questions that are suited to be examined by computational methods with big data
5. Write a research article that utilizes the techniques and methods of computational social science to address social science problems, or design a project that use computational social science to address some real-world problems.

Study Goals

1. Describe the opportunities and challenges of computational social science
2. Evaluate computational social science research on social phenomena
3. Practice the essential techniques to analyze social big data, especially text data (covered in Tutorials)
 - Getting data
 - Managing data
4. Propose research questions that are suited to be examined by computational methods with big data
5. Write a research article that utilizes the techniques and methods of computational social science to address social science problems, or design a project that use computational social science to address some real-world problems.

Study Goals

1. Describe the opportunities and challenges of computational social science
2. Evaluate computational social science research on social phenomena
3. Practice the essential techniques to analyze social big data, especially text data (covered in Tutorials)
 - Getting data
 - Managing data
 - Analyzing data with appropriate methods
4. Propose research questions that are suited to be examined by computational methods with big data
5. Write a research article that utilizes the techniques and methods of computational social science to address social science problems, or design a project that use computational social science to address some real-world problems.

Instructors

Instructor: ZHANG, Han

- Office: 2379

Teaching Assistant: Chen, Pei

Instructors

Instructor: ZHANG, Han

- Office: 2379
- Email: zhangh@ust.hk

Teaching Assistant: Chen, Pei

Instructors

Instructor: ZHANG, Han

- Office: 2379
- Email: zhangh@ust.hk
- Office Hour: Mon 11-12AM (or email me to find a time)

Teaching Assistant: Chen, Pei

Instructors

Instructor: ZHANG, Han

- Office: 2379
- Email: zhangh@ust.hk
- Office Hour: Mon 11-12AM (or email me to find a time)
 - In-office

Teaching Assistant: Chen, Pei

Instructors

Instructor: ZHANG, Han

- Office: 2379
- Email: zhangh@ust.hk
- Office Hour: Mon 11-12AM (or email me to find a time)
 - In-office
 - Zoom: <https://hkust.zoom.us/j/6522716568>

Teaching Assistant: Chen, Pei

Instructors

Instructor: ZHANG, Han

- Office: 2379
- Email: zhangh@ust.hk
- Office Hour: Mon 11-12AM (or email me to find a time)
 - In-office
 - Zoom: <https://hkust.zoom.us/j/6522716568>

Teaching Assistant: Chen, Pei

- Office: 3001

Instructors

Instructor: ZHANG, Han

- Office: 2379
- Email: zhangh@ust.hk
- Office Hour: Mon 11-12AM (or email me to find a time)
 - In-office
 - Zoom: <https://hkust.zoom.us/j/6522716568>

Teaching Assistant: Chen, Pei

- Office: 3001
- Email: pchenam@connect.ust.hk

Instructors

Instructor: ZHANG, Han

- Office: 2379
- Email: zhangh@ust.hk
- Office Hour: Mon 11-12AM (or email me to find a time)
 - In-office
 - Zoom: <https://hkust.zoom.us/j/6522716568>

Teaching Assistant: Chen, Pei

- Office: 3001
- Email: pchenam@connect.ust.hk
- Office Hour: TBD

Online teaching

- Some rules:

Online teaching

- Some rules:
 - Keep the video on (this will be counted as class participation)

Online teaching

- Some rules:
 - Keep the video on (this will be counted as class participation)
 - Mute yourself while not speaking

Online teaching

- Some rules:
 - Keep the video on (this will be counted as class participation)
 - Mute yourself while not speaking
 - Feel free to stop me at any time if you have any question

Online teaching

- Some rules:
 - Keep the video on (this will be counted as class participation)
 - Mute yourself while not speaking
 - Feel free to stop me at any time if you have any question
 - Recordings will be available after class

Course material

- We will use GitHub for teaching

Course material

- We will use GitHub for teaching
- Lecture material and tutorial will be available at:
<https://github.com/HKUST-SOSC4300-5500>

Course material

- We will use GitHub for teaching
- Lecture material and tutorial will be available at:
<https://github.com/HKUST-SOSC4300-5500>
- Please use the version on GitHub as the authoritative version

Schedule (tentative)

Week	Topic
1	02-08 Introduction; big data
2	02-15 Prediction; Machine learning
3	02-22 Prediction; Evaluation
4	03-01 Text (I)
5	03-08 Text (II); supervised
6	03-15 Text (III); embedding
7	03-22 Text (IV); unsupervised
8	03-29 Image data
9	04-12 Network; agent-based modeling
10	04-19 Network; strength of ties
11	04-26 Network; causal inference

Grading Components

- 4300 and 5500 will be graded independently

Attendance and participation	10%
Homework assignments	30%
Literature review Report	15%
Final Paper/Project	
Presentation	15%
Write-up	30%

Grouping

- You should finish all tasks in groups

Grouping

- You should finish all tasks in groups
 - If there is any **MPhil or PhD** student in a group: max group size is 2

Grouping

- You should finish all tasks in groups
 - If there is any **MPhil or PhD** student in a group: max group size is 2
 - Otherwise: 3 to 4 in a group (e.g., 4 UG in a group)

Grouping

- You should finish all tasks in groups
 - If there is any **MPhil or PhD** student in a group: max group size is 2
 - Otherwise: 3 to 4 in a group (e.g., 4 UG in a group)
- Finish grouping by **the end of February**

Grouping

- You should finish all tasks in groups
 - If there is any **MPhil or PhD** student in a group: max group size is 2
 - Otherwise: 3 to 4 in a group (e.g., 4 UG in a group)
- Finish grouping by **the end of February**
 - we will have first assignment then

Class participation

- Read the required readings on syllabus before each week's lecture

Class participation

- Read the required readings on syllabus before each week's lecture
- Answer questions about the assigned readings

Class participation

- Read the required readings on syllabus before each week's lecture
- Answer questions about the assigned readings
- Ask questions about the parts you did not understand.

Class participation

- Read the required readings on syllabus before each week's lecture
- Answer questions about the assigned readings
- Ask questions about the parts you did not understand.
- If you are uncomfortable speaking up in class, send the question in Zoom's chat window, post them on Github, come to my office hours, or send your questions to instructors via e-mail.

Homework assignments

- There will be 3 to 4 coding exercise as homework.

Homework assignments

- There will be 3 to 4 coding exercise as homework.
- These homework assignments test your knowledge of analyzing data using statistical software.

Homework assignments

- There will be 3 to 4 coding exercise as homework.
- These homework assignments test your knowledge of analyzing data using statistical software.
- Each exercise is due in **two weeks** after the release of assignment.

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as “social media texts and sentiment analysis”

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as "social media texts and sentiment analysis"
- Some examples of research areas:

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as "social media texts and sentiment analysis"
- Some examples of research areas:
 - Sociology: internal or international migration, social inequality, race and ethnicity relations, wellbeing

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as “social media texts and sentiment analysis”
- Some examples of research areas:
 - Sociology: internal or international migration, social inequality, race and ethnicity relations, wellbeing
 - Political science: government performance, government policy effectiveness, election, protests and social movements

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as "social media texts and sentiment analysis"
- Some examples of research areas:
 - Sociology: internal or international migration, social inequality, race and ethnicity relations, wellbeing
 - Political science: government performance, government policy effectiveness, election, protests and social movements
 - Economics: measuring economic growth with big data

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as “social media texts and sentiment analysis”
- Some examples of research areas:
 - Sociology: internal or international migration, social inequality, race and ethnicity relations, wellbeing
 - Political science: government performance, government policy effectiveness, election, protests and social movements
 - Economics: measuring economic growth with big data
 - History: historical development of an idea

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as "social media texts and sentiment analysis"
- Some examples of research areas:
 - Sociology: internal or international migration, social inequality, race and ethnicity relations, wellbeing
 - Political science: government performance, government policy effectiveness, election, protests and social movements
 - Economics: measuring economic growth with big data
 - History: historical development of an idea
 - Psychology: measuring personality with big data

Literature Review

- Select a social topic and summarize how researchers have used computational methods and/or big data to study this particular research area.
- **Be specific:** don't choose topics such as “social media texts and sentiment analysis”
- Some examples of research areas:
 - Sociology: internal or international migration, social inequality, race and ethnicity relations, wellbeing
 - Political science: government performance, government policy effectiveness, election, protests and social movements
 - Economics: measuring economic growth with big data
 - History: historical development of an idea
 - Psychology: measuring personality with big data
 - Communication and information science: content and spread of fake news/hate speeches on social media

Literature Review

- **Written Report:** each literature review report should contain at least **8 pages, 12 points, double space.**

Literature Review

- **Written Report:** each literature review report should contain at least **8 pages, 12 points, double space.**
- Use this to prepare your final paper

Final Paper/Project

- You can choose to write a research final paper

Final Paper/Project

- You can choose to write a research final paper
 - The intended audience for research final paper are other **researchers**

Final Paper/Project

- You can choose to write a research final paper
 - The intended audience for research final paper are other **researchers**
- Or a project that analyze a “real-world” case.

Final Paper/Project

- You can choose to write a research final paper
 - The intended audience for research final paper are other **researchers**
- Or a project that analyze a “real-world” case.
 - Like building a website that has nice visualization

Final Paper/Project

- You can choose to write a research final paper
 - The intended audience for research final paper are other **researchers**
- Or a project that analyze a “real-world” case.
 - Like building a website that has nice visualization
 - Project should attract **layman**

Final Paper/Project

- If you choose to write an research paper:

Final Paper/Project

- If you choose to write a research paper:
 - Presentation (20 minutes): follow a standard presentation style for academic talks.

Final Paper/Project

- If you choose to write an research paper:
 - Presentation (20 minutes): follow a standard presentation style for academic talks.
 - Final paper/project: 20 pages, 12 points, double space, including Tables, Figures and References.

Final Paper/Project

- If you choose to write an research paper:
 - Presentation (20 minutes): follow a standard presentation style for academic talks.
 - Final paper/project: 20 pages, 12 points, double space, including Tables, Figures and References.
- If you choose to do a project:

Final Paper/Project

- If you choose to write an research paper:
 - Presentation (20 minutes): follow a standard presentation style for academic talks.
 - Final paper/project: 20 pages, 12 points, double space, including Tables, Figures and References.
- If you choose to do a project:
 - Presentation (20 minutes). Show case your project in front of the class.

Final Paper/Project

- If you choose to write an research paper:
 - Presentation (20 minutes): follow a standard presentation style for academic talks.
 - Final paper/project: 20 pages, 12 points, double space, including Tables, Figures and References.
- If you choose to do a project:
 - Presentation (20 minutes). Show case your project in front of the class.
 - Technical report: a short write up on short background/dataset/methods; 10 pages, 12 points, double space;

Final Paper/Project

- If you choose to write an research paper:
 - Presentation (20 minutes): follow a standard presentation style for academic talks.
 - Final paper/project: 20 pages, 12 points, double space, including Tables, Figures and References.
- If you choose to do a project:
 - Presentation (20 minutes). Show case your project in front of the class.
 - Technical report: a short write up on short background/dataset/methods; 10 pages, 12 points, double space;
- If there are so many groups, we will possibly shift to poster

We will use GitHub for teaching

- We will use GitHub for teaching

We will use GitHub for teaching

- We will use GitHub for teaching
- GitHub is commonly used among researchers and companies

We will use GitHub for teaching

- We will use GitHub for teaching
- GitHub is commonly used among researchers and companies
- It allows us to manage individual research workflow more easily

We will use GitHub for teaching

- We will use GitHub for teaching
- GitHub is commonly used among researchers and companies
- It allows us to manage individual research workflow more easily
- And allows team work easily (especially if you have lots of codes/data)

We will use GitHub for teaching

- We will use GitHub for teaching
- GitHub is commonly used among researchers and companies
- It allows us to manage individual research workflow more easily
- And allows team work easily (especially if you have lots of codes/data)
- Last, GitHub simplifies sharing data and codes, making research more transparent and useful for the community

So what is it?

Linus Torvalds: creator of Linux



When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control
 - We will update a file, but want to keep some old version if we need to roll back

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control
 - We will update a file, but want to keep some old version if we need to roll back
 - Renaming files to file.r, file_ver1.r, file_ver2.r ...: an ugly solution

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control
 - We will update a file, but want to keep some old version if we need to roll back
 - Renaming files to file.r, file_ver1.r, file_ver2.r ...: an ugly solution
- Project management:

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control
 - We will update a file, but want to keep some old version if we need to roll back
 - Renaming files to file.r, file_ver1.r, file_ver2.r ...: an ugly solution
- Project management:
 - Multiple people working on the same file

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control
 - We will update a file, but want to keep some old version if we need to roll back
 - Renaming files to file.r, file_ver1.r, file_ver2.r ...: an ugly solution
- Project management:
 - Multiple people working on the same file
 - Contributors send Linus code chunks, and Linus put these code chunks to its appropriate place manually

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control
 - We will update a file, but want to keep some old version if we need to roll back
 - Renaming files to file.r, file_ver1.r, file_ver2.r ...: an ugly solution
- Project management:
 - Multiple people working on the same file
 - Contributors send Linus code chunks, and Linus put these code chunks to its appropriate place manually
 - Error prone: over 3,500 lines of codes added per day

When Linux grows

- Over 30,000 files; 15 Million lines of codes in 2012
- Version control
 - We will update a file, but want to keep some old version if we need to roll back
 - Renaming files to file.r, file_ver1.r, file_ver2.r ...: an ugly solution
- Project management:
 - Multiple people working on the same file
 - Contributors send Linux code chunks, and Linus put these code chunks to its appropriate place manually
 - Error prone: over 3,500 lines of codes added per day
- So Linus intended **Git**

GitHub

- Git is free;

GitHub

- Git is free;
- GitHub is an online service that stores Git repository;

GitHub

- Git is free;
- GitHub is an online service that stores Git repository;
 - other options: GitLab.

GitHub

- Git is free;
- GitHub is an online service that stores Git repository;
 - other options: GitLab.
 - Student can get a free version, allowing you to obtain **private** repo.

GitHub

- Git is free;
- GitHub is an online service that stores Git repository;
 - other options: GitLab.
 - Student can get a free version, allowing you to obtain **private** repo.
- GitHub also allows us to share codes and data easily

GitHub

- Git is free;
- GitHub is an online service that stores Git repository;
 - other options: GitLab.
 - Student can get a free version, allowing you to obtain **private** repo.
- GitHub also allows us to share codes and data easily
- `https://github.com/HKUST-SOSC4300-5500/Tutorial-Material/blob/master/week1/1-Hello-World.ipynb`

Basic operations

- **add** : start to track changes of some files

Basic operations

- **add** : start to track changes of some files
- **commit** : record changes **locally**

Basic operations

- **add** : start to track changes of some files
- **commit** : record changes **locally**
 - First commit records the oldest version of a file

Basic operations

- **add** : start to track changes of some files
- **commit** : record changes **locally**
 - First commit records the oldest version of a file
 - Newest commit records the most recent version of a file

Basic operations

- **add** : start to track changes of some files
- **commit** : record changes **locally**
 - First commit records the oldest version of a file
 - Newest commit records the most recent version of a file
 - Only changes in a file will be saved; saving disk spaces

Basic operations

- **add** : start to track changes of some files
- **commit** : record changes **locally**
 - First commit records the oldest version of a file
 - Newest commit records the most recent version of a file
 - Only changes in a file will be saved; saving disk spaces
- **push** : send committed changes to GitHub so that others can view it

Basic operations

- **add** : start to track changes of some files
- **commit** : record changes **locally**
 - First commit records the oldest version of a file
 - Newest commit records the most recent version of a file
 - Only changes in a file will be saved; saving disk spaces
- **push** : send committed changes to GitHub so that others can view it
- **pull** : sync back changes from remote to your local machine

Git for version control: example

- Say you are working on assignment 1, with a R script `homework1.r`

Git for version control: example

- Say you are working on assignment 1, with a R script `homework1.r`
- Step 1: **add** `homework1.r` to tell git to pay attention to changes in this file

Git for version control: example

- Say you are working on assignment 1, with a R script `homework1.r`
- Step 1: **add** `homework1.r` to tell git to pay attention to changes in this file
- Whenever you want to save a version for backup, **commit** `homework1.r`

Git for version control: example

- Say you are working on assignment 1, with a R script `homework1.r`
- Step 1: **add** `homework1.r` to tell git to pay attention to changes in this file
- Whenever you want to save a version for backup, **commit** `homework1.r`
- When you feel you are ready to submit, **push** to GitHub

Git for version control: example

- Say you are working on assignment 1, with a R script `homework1.r`
- Step 1: **add** `homework1.r` to tell git to pay attention to changes in this file
- Whenever you want to save a version for backup, **commit** `homework1.r`
- When you feel you are ready to submit, **push** to GitHub
 - Of course you can push multiple times, if deadline is not passed.

Git for project management

- Project management for groups: (called repository or repo)

which allow the owner to *pull* changes from local branch to master branch

Git for project management

- Project management for groups: (called repository or repo)
 - Owner controls **master branch**; this is the central and authoritative version

which allow the owner to *pull* changes from local branch to master branch

Git for project management

- Project management for groups: (called repository or repo)
 - Owner controls **master branch**; this is the central and authoritative version
 - Every member works with his own local version (called **local branch**)

which allow the owner to *pull* changes from local branch to master branch

Git for project management

- Project management for groups: (called repository or repo)
 - Owner controls **master branch**; this is the central and authoritative version
 - Every member works with his own local version (called **local branch**)
 - Only owner can directly **push** to the central version

which allow the owner to *pull* changes from local branch to master branch

Git for project management

- Project management for groups: (called repository or repo)
 - Owner controls **master branch**; this is the central and authoritative version
 - Every member works with his own local version (called **local branch**)
 - Only owner can directly **push** to the central version
 - Members can create **pull request**,

which allow the owner to *pull* changes from local branch to master branch

Why Git, why not other file services, for version control?

- Say, why not Dropbox or Google Drive, which also keeps track of file versions?

Why Git, why not other file services, for version control?

- Say, why not Dropbox or Google Drive, which also keeps track of file versions?
- Scenario 1: I only want to let others see my changes when I am ready to.

Why Git, why not other file services, for version control?

- Say, why not Dropbox or Google Drive, which also keeps track of file versions?
- Scenario 1: I only want to let others see my changes when I am ready to.
- Scenario 2: There is a clear hierarchy in the project, with someone can accept or reject your changes.

Three ways to use Git

- Command line: for advanced and interested users

Three ways to use Git

- Command line: for advanced and interested users
- GitHub's official app. MacOS and Windows version can be downloaded at <https://docs.github.com/en/desktop/installing-and-configuring-github-desktop/installing-github-desktop>

Three ways to use Git

- Command line: for advanced and interested users
- GitHub's official app. MacOS and Windows version can be downloaded at <https://docs.github.com/en/desktop/installing-and-configuring-github-desktop/installing-github-desktop>
 - Recommended

Three ways to use Git

- Command line: for advanced and interested users
- GitHub's official app. MacOS and Windows version can be downloaded at <https://docs.github.com/en/desktop/installing-and-configuring-github-desktop/installing-github-desktop>
 - **Recommended**
- Directly change and upload files on GitHub website

Three ways to use Git

- Command line: for advanced and interested users
- GitHub's official app. MacOS and Windows version can be downloaded at <https://docs.github.com/en/desktop/installing-and-configuring-github-desktop/installing-github-desktop>
 - **Recommended**
- Directly change and upload files on GitHub website
 - Note: if you directly upload a file on GitHub in your own repo, you are performing add, commit and push operation simultaneously

How we use GitHub for this class

- Lecture and tutorial material will be available at <https://github.com/HKUST-SOSC4300-5500>

How we use GitHub for this class

- Lecture and tutorial material will be available at <https://github.com/HKUST-SOSC4300-5500>
 - You can view and download the files

How we use GitHub for this class

- Lecture and tutorial material will be available at <https://github.com/HKUST-SOSC4300-5500>
 - You can view and download the files
- Assignment will be submitted on GitHub:
<https://classroom.github.com/classrooms/70853257-hkust-sosc4300-5500-classroom-1>

How we use GitHub for this class

- Lecture and tutorial material will be available at <https://github.com/HKUST-SOSC4300-5500>
 - You can view and download the files
- Assignment will be submitted on GitHub:
<https://classroom.github.com/classrooms/70853257-hkust-sosc4300-5500-classroom-1>
 - Homework will be **private**: only you and instructors can see your commits

How we use GitHub for this class

- Lecture and tutorial material will be available at <https://github.com/HKUST-SOSC4300-5500>
 - You can view and download the files
- Assignment will be submitted on GitHub:
<https://classroom.github.com/classrooms/70853257-hkust-sosc4300-5500-classroom-1>
 - Homework will be **private**: only you and instructors can see your commits
 - Literature review and final paper/project will be **public**: since these knowledge will be beneficial to others

How we use GitHub for this class

- Lecture and tutorial material will be available at <https://github.com/HKUST-SOSC4300-5500>
 - You can view and download the files
- Assignment will be submitted on GitHub:
<https://classroom.github.com/classrooms/70853257-hkust-sosc4300-5500-classroom-1>
 - Homework will be **private**: only you and instructors can see your commits
 - Literature review and final paper/project will be **public**: since these knowledge will be beneficial to others
 - We will set deadlines; your last push before the deadline will be automatically treated as your final submitted version

Data acquisition

- Access provided by Government/Company

Data acquisition

- Access provided by Government/Company
- Collect by yourself

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface
 - Easy; websites will provide you instructions to obtain some of their data

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface
 - Easy; websites will provide you instructions to obtain some of their data
 - Restricted by what the data provider allows you to download

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface
 - Easy; websites will provide you instructions to obtain some of their data
 - Restricted by what the data provider allows you to download
 - Web crawling

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface
 - Easy; websites will provide you instructions to obtain some of their data
 - Restricted by what the data provider allows you to download
 - Web crawling
 - In principle, you can download whatever you can see as a human

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface
 - Easy; websites will provide you instructions to obtain some of their data
 - Restricted by what the data provider allows you to download
 - Web crawling
 - In principle, you can download whatever you can see as a human
 - Technically, much more challenging

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface
 - Easy; websites will provide you instructions to obtain some of their data
 - Restricted by what the data provider allows you to download
 - Web crawling
 - In principle, you can download whatever you can see as a human
 - Technically, much more challenging
 - Only use this when there is no API / API does not satisfy your need

Data acquisition

- Access provided by Government/Company
- Collect by yourself
 - Manual download
 - only works if you are dealing with small data
 - API: Application Programming Interface
 - Easy; websites will provide you instructions to obtain some of their data
 - Restricted by what the data provider allows you to download
 - Web crawling
 - In principle, you can download whatever you can see as a human
 - Technically, much more challenging
 - Only use this when there is no API / API does not satisfy your need
- We will cover these in tutorials

Some dataset archives

- Some collection of datasets:

Some dataset archives

- Some collection of datasets:
 - Google Dataset search:
<https://datasetsearch.research.google.com/>

Some dataset archives

- Some collection of datasets:
 - Google Dataset search:
<https://datasetsearch.research.google.com/>
 - Kaggle Dataset: <https://www.kaggle.com/datasets?tags=14104-text+data>

Some dataset archives

- Some collection of datasets:
 - Google Dataset search:
<https://datasetsearch.research.google.com/>
 - Kaggle Dataset: <https://www.kaggle.com/datasets?tags=14104-text+data>
 - UCI's machine learning repository:
<https://archive.ics.uci.edu/ml/datasets.php>

Some dataset archives

- Some collection of datasets:
 - Google Dataset search:
<https://datasetsearch.research.google.com/>
 - Kaggle Dataset: <https://www.kaggle.com/datasets?tags=14104-text+data>
 - UCI's machine learning repository:
<https://archive.ics.uci.edu/ml/datasets.php>
 - US patents: <https://www.google.com/googlebooks/uspto-patents-assignments.html>

Some dataset archives

- Some collection of datasets:
 - Google Dataset search:
<https://datasetsearch.research.google.com/>
 - Kaggle Dataset: <https://www.kaggle.com/datasets?tags=14104-text+data>
 - UCI's machine learning repository:
<https://archive.ics.uci.edu/ml/datasets.php>
 - US patents: <https://www.google.com/googlebooks/uspto-patents-assignments.html>
 - Wikipedia texts:
https://en.wikipedia.org/wiki/Wikipedia:Database_download

Roadmap

- We will discuss 10 characteristics of big data, following

Roadmap

- We will discuss 10 characteristics of big data, following
- Chapter 2, Matthew Salganik, *Bit by Bit: Social Research in the Digital Age*, Princeton University Press, 2019

Roadmap

- We will discuss 10 characteristics of big data, following
- Chapter 2, Matthew Salganik, *Bit by Bit: Social Research in the Digital Age*, Princeton University Press, 2019
- <https://www.bitbybitbook.com/en/1st-ed/observing-behavior/data/>

Roadmap

- We will discuss 10 characteristics of big data, following
- Chapter 2, Matthew Salganik, *Bit by Bit: Social Research in the Digital Age*, Princeton University Press, 2019
- <https://www.bitbybitbook.com/en/1st-ed/observing-behavior/data/>
- After our discussions, you can critically evaluate pros and cons of big data

Characteristics 1: Big

- Jean-Baptiste Michel, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K. Gray, The Google Books Team, Joseph P. Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, Jon Orwant, Steven Pinker, Martin A. Nowak, and Erez Lieberman Aiden, *Quantitative Analysis of Culture Using Millions of Digitized Books*, Science **331** (2011), no. 6014, 176–182

[Our] corpus contains over 500 billion words, in English (361 billion), French (45 billion), Spanish (45 billion), German (37 billion), Chinese (13 billion), Russian (35 billion), and Hebrew (2 billion). The oldest works were published in the 1500s. The early decades are represented by only a few books per year, comprising several hundred thousand words. By 1800, the corpus grows to 98 million words per year; by 1900, 1.8 billion; and by 2000, 11 billion. The corpus cannot be read by a human. If you tried to read only English-language entries from the year 2000 alone, at the reasonable pace of 200 words/min, without interruptions for food or sleep, it would take 80 years. The sequence of letters is 1000 times longer than the human genome: If you wrote it out in a straight line, it would reach to the Moon and back 10 times over.

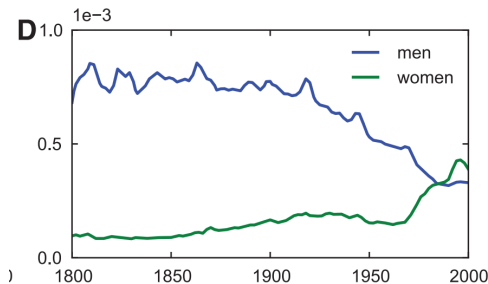
Characteristics 1: Big

- Jean-Baptiste Michel, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K. Gray, The Google Books Team, Joseph P. Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, Jon Orwant, Steven Pinker, Martin A. Nowak, and Erez Lieberman Aiden, *Quantitative Analysis of Culture Using Millions of Digitized Books*, Science **331** (2011), no. 6014, 176–182
- They turned Google Books into word counts and released the data

[Our] corpus contains over 500 billion words, in English (361 billion), French (45 billion), Spanish (45 billion), German (37 billion), Chinese (13 billion), Russian (35 billion), and Hebrew (2 billion). The oldest works were published in the 1500s. The early decades are represented by only a few books per year, comprising several hundred thousand words. By 1800, the corpus grows to 98 million words per year; by 1900, 1.8 billion; and by 2000, 11 billion. The corpus cannot be read by a human. If you tried to read only English-language entries from the year 2000 alone, at the reasonable pace of 200 words/min, without interruptions for food or sleep, it would take 80 years. The sequence of letters is 1000 times longer than the human genome: If you wrote it out in a straight line, it would reach to the Moon and back 10 times over.

Characteristics 1: Big

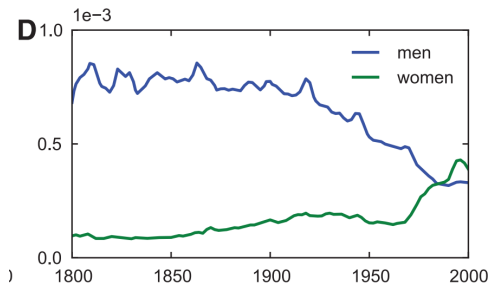
- Explore their project here:
<https://books.google.com/ngrams>



- [In class discussion]: do we really need this many data to draw the conclusion that women's right are rising? Can't we use smaller data to reach the same conclusion?

Characteristics 1: Big

- Explore their project here:
<https://books.google.com/ngrams>
- E.g., “In the battle of the sexes, the women are gaining ground on the men”



- [In class discussion]: do we really need this many data to draw the conclusion that women's right are rising? Can't we use smaller data to reach the same conclusion?

Characteristics 1: Big

- Big data are good at showing heterogeneity, which cannot be obtained by smaller data



Characteristics 1: Big

- Big data are good at showing heterogeneity, which cannot be obtained by smaller data
- (Chetty et al. 2014), estimates of a child's chances of reaching the top 20% of income distribution given parents in the bottom 20% .



Characteristics 1: Big

- Big data are good at showing heterogeneity, which cannot be obtained by smaller data
- (Chetty et al. 2014), estimates of a child's chances of reaching the top 20% of income distribution given parents in the bottom 20% .
- “The regional-level estimates, which show heterogeneity, naturally lead to interesting and important questions that do not arise from a single national-level estimate. These regional-level estimates were made possible in part because the researchers were using a large big data source: the tax records of 40 million people.



Characteristics 2: Always-on

- Traditional data survey: once a year, or on demand

Characteristics 2: Always-on

- Traditional data survey: once a year, or on demand
- Big data: always-on measure

Characteristics 2: Always-on

- Traditional data survey: once a year, or on demand
- Big data: always-on measure
- Ceren Budak and Duncan Watts, *Dissecting the Spirit of Gezi: Influence vs. Selection in the Occupy Gezi Movement*, Sociological Science **2** (2015), 370–397

Characteristics 2: Always-on

- Traditional data survey: once a year, or on demand
- Big data: always-on measure
- Ceren Budak and Duncan Watts, *Dissecting the Spirit of Gezi: Influence vs. Selection in the Occupy Gezi Movement*, Sociological Science **2** (2015), 370–397
- What kinds of people were more likely to participate in the Gezi protests in 2013?

Characteristics 2: Always-on

- Traditional data survey: once a year, or on demand
- Big data: always-on measure
- Ceren Budak and Duncan Watts, *Dissecting the Spirit of Gezi: Influence vs. Selection in the Occupy Gezi Movement*, Sociological Science **2** (2015), 370–397
- What kinds of people were more likely to participate in the Gezi protests in 2013?
- Whether participation changed attitudes of participants and nonparticipants differently?

Characteristics 2: Always-on

- Traditional data survey: once a year, or on demand
- Big data: always-on measure
- Ceren Budak and Duncan Watts, *Dissecting the Spirit of Gezi: Influence vs. Selection in the Occupy Gezi Movement*, Sociological Science **2** (2015), 370–397
- What kinds of people were more likely to participate in the Gezi protests in 2013?
- Whether participation changed attitudes of participants and nonparticipants differently?
- Hard with survey data:

Characteristics 2: Always-on

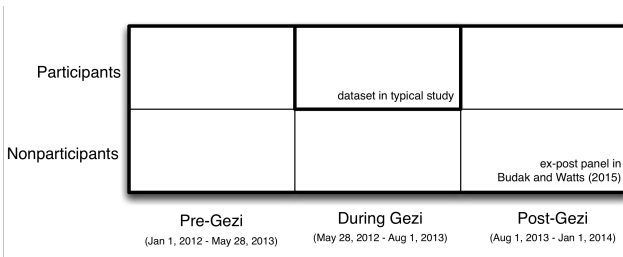
- Traditional data survey: once a year, or on demand
- Big data: always-on measure
- Ceren Budak and Duncan Watts, *Dissecting the Spirit of Gezi: Influence vs. Selection in the Occupy Gezi Movement*, Sociological Science **2** (2015), 370–397
- What kinds of people were more likely to participate in the Gezi protests in 2013?
- Whether participation changed attitudes of participants and nonparticipants differently?
- Hard with survey data:
 - You cannot predict when a protest occur, and thus cannot get **pre-protest** information

Characteristics 2: Always-on

- Traditional data survey: once a year, or on demand
- Big data: always-on measure
- Ceren Budak and Duncan Watts, *Dissecting the Spirit of Gezi: Influence vs. Selection in the Occupy Gezi Movement*, Sociological Science **2** (2015), 370–397
- What kinds of people were more likely to participate in the Gezi protests in 2013?
- Whether participation changed attitudes of participants and nonparticipants differently?
- Hard with survey data:
 - You cannot predict when a protest occur, and thus cannot get **pre-protest** information
 - It's also not easy to get samples of non-participants

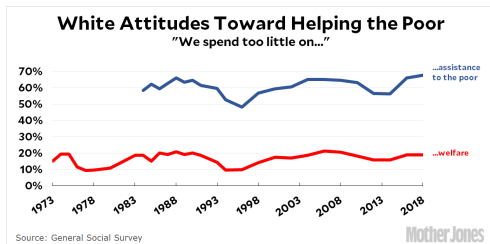
Characteristics 2: Always-On

- Using geolocated posts on Twitter



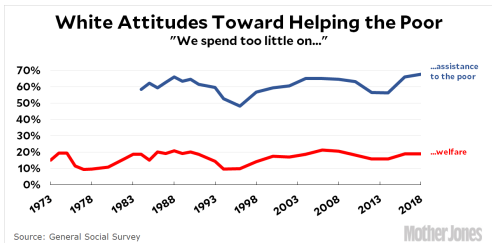
Characteristics 3: Non-Reactive

- Big data are mostly obtained **unobtrusively**;



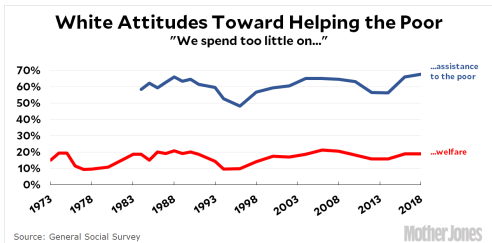
Characteristics 3: Non-Reactive

- Big data are mostly obtained **unobtrusively**;
 - People are generally not aware that their data are being captured



Characteristics 3: Non-Reactive

- Big data are mostly obtained **unobtrusively**;
 - People are generally not aware that their data are being captured
- Survey and lab experiments obtain data **obtrusively**, and results often depend on how you ask



Characteristics 3: Non-Reactive

- [In class activity]: is non-reactiveness always good?
- What people put on social media may be just showing off, not their daily lives

Characteristics 3: Non-Reactive

- [In class activity]: is non-reactiveness always good?
- What people put on social media may be just showing off, not their daily lives
- Sometimes it is quicker to ask, especially for questions that are less likely to vary depending on how you ask

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want
- Three types of data are especially likely to be missing

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want
- Three types of data are especially likely to be missing
 - demographic information

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want
- Three types of data are especially likely to be missing
 - demographic information
 - E.g., Google N-grams has gigantic dataset, but does not directly has author’s biography

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want
- Three types of data are especially likely to be missing
 - demographic information
 - E.g., Google N-grams has gigantic dataset, but does not directly has author’s biography
 - behavior on other platforms

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want
- Three types of data are especially likely to be missing
 - demographic information
 - E.g., Google N-grams has gigantic dataset, but does not directly has author’s biography
 - behavior on other platforms
 - data to operationalize theoretical concepts

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want
- Three types of data are especially likely to be missing
 - demographic information
 - E.g., Google N-grams has gigantic dataset, but does not directly has author’s biography
 - behavior on other platforms
 - data to operationalize theoretical concepts
 - people who are more intelligent earn more money

Characteristics 4: Incomplete

- “No matter how big your big data, it probably doesn’t have the information you want”
- This is a property of **re-purposing** the data;
 - for survey and lab experiments, you can in principle ask what you want
- Three types of data are especially likely to be missing
 - demographic information
 - E.g., Google N-grams has gigantic dataset, but does not directly has author’s biography
 - behavior on other platforms
 - data to operationalize theoretical concepts
 - people who are more intelligent earn more money
 - how do you measure intelligence with big data? Not easy

Characteristics 4: Incomplete

- But incompleteness may also occur for traditional survey data.
- [in class activity] E.g., a classical argument in social networks: the more centered you are in social network, the more wealthy you are

Characteristics 4: Incomplete

- But incompleteness may also occur for traditional survey data.
 - A national representative survey does not
- [in class activity] E.g., a classical argument in social networks: the more centered you are in social network, the more wealthy you are

Characteristics 4: Incomplete

- But incompleteness may also occur for traditional survey data.
 - A national representative survey does not
- [in class activity] E.g., a classical argument in social networks: the more centered you are in social network, the more wealthy you are
 - How do we test this argument with survey data?

Characteristics 4: Incomplete

- But incompleteness may also occur for traditional survey data.
 - A national representative survey does not
- [in class activity] E.g., a classical argument in social networks: the more centered you are in social network, the more wealthy you are
 - How do we test this argument with survey data?
 - How can we test this argument with big data?

Characteristics 4: Incomplete

- But incompleteness may also occur for traditional survey data.
 - A national representative survey does not
- [in class activity] E.g., a classical argument in social networks: the more centered you are in social network, the more wealthy you are
 - How do we test this argument with survey data?
 - How can we test this argument with big data?
 - What will be the best data source you can think of?

Characteristics 5: Inaccessible

- Many useful data are not directly available to researchers; they are stored in government and company servers

Characteristics 5: Inaccessible

- Many useful data are not directly available to researchers; they are stored in government and company servers
- Reason 1: commercial/government secrets

Characteristics 5: Inaccessible

- Many useful data are not directly available to researchers; they are stored in government and company servers
- Reason 1: commercial/government secrets
- Reason 2: terms-of-service agreements

Characteristics 5: Inaccessible

- Many useful data are not directly available to researchers; they are stored in government and company servers
- Reason 1: commercial/government secrets
- Reason 2: terms-of-service agreements
- Reason 3: releasing data sometimes lead to privacy concerns

Characteristics 6: Nonrepresentative

- Many big data sources are not representative samples from some well-defined population
- [in class activity] So does it mean big data are not useful?
When nonrepresentative data are useful?

Characteristics 7: Drifting

- Digital world changes so quick so that it's still too early to use big data to study long-term trends

Characteristics 7: Drifting

- Digital world changes so quick so that it's still too early to use big data to study long-term trends
- Population drift (change in who is using them)

Characteristics 7: Drifting

- Digital world changes so quick so that it's still too early to use big data to study long-term trends
- Population drift (change in who is using them)
- Behavioral drift (change in how people are using them)

Characteristics 7: Drifting

- Digital world changes so quick so that it's still too early to use big data to study long-term trends
- Population drift (change in who is using them)
- Behavioral drift (change in how people are using them)
- System drift (change in the system itself).

Characteristics 8: Algorithm Confounded

- Again, gov/companies control the data generating process

Characteristics 8: Algorithm Confounded

- Again, gov/companies control the data generating process
- E.g., to what extent your friends of friends are more likely to be your friends?

Characteristics 8: Algorithm Confounded

- Again, gov/companies control the data generating process
- E.g., to what extent your friends of friends are more likely to be your friends?
 - It's called *triad closure* in social network analysis

Characteristics 8: Algorithm Confounded

- Again, gov/companies control the data generating process
- E.g., to what extent your friends of friends are more likely to be your friends?
 - It's called *triad closure* in social network analysis
- With social networks such as Facebook, it's possible to empirically measure this quantity precisely

Characteristics 8: Algorithm Confounded

- Again, gov/companies control the data generating process
- E.g., to what extent your friends of friends are more likely to be your friends?
 - It's called *triad closure* in social network analysis
- With social networks such as Facebook, it's possible to empirically measure this quantity precisely
- Until Facebook began to recommend friends to users

Characteristics 8: Algorithm Confounded

- Again, gov/companies control the data generating process
- E.g., to what extent your friends of friends are more likely to be your friends?
 - It's called *triad closure* in social network analysis
- With social networks such as Facebook, it's possible to empirically measure this quantity precisely
- Until Facebook began to recommend friends to users
- Now, are what we observe because of Facebook's recommendation or innate tendency for friends of friends to become friends?

Characteristics 9: Dirty

- “Big data sources can be loaded with junk and spam”

Characteristics 9: Dirty

- “Big data sources can be loaded with junk and spam”
- Example: Eitan Adam Pechenick, Christopher M. Danforth, and Peter Sheridan Dodds, *Characterizing the Google Books Corpus: Strong Limits to Inferences of Socio-Cultural and Linguistic Evolution*, PLOS ONE **10** (2015), no. 10, e0137041

Characteristics 9: Dirty

- “Big data sources can be loaded with junk and spam”
- Example: Eitan Adam Pechenick, Christopher M. Danforth, and Peter Sheridan Dodds, *Characterizing the Google Books Corpus: Strong Limits to Inferences of Socio-Cultural and Linguistic Evolution*, PLOS ONE **10** (2015), no. 10, e0137041
- Problem 1: OCR error

Characteristics 9: Dirty

- “Big data sources can be loaded with junk and spam”
- Example: Eitan Adam Pechenick, Christopher M. Danforth, and Peter Sheridan Dodds, *Characterizing the Google Books Corpus: Strong Limits to Inferences of Socio-Cultural and Linguistic Evolution*, PLOS ONE **10** (2015), no. 10, e0137041
- Problem 1: OCR error
- The count of F-word between 1800 to 2000

Characteristics 9: Dirty

- “Big data sources can be loaded with junk and spam”
- Example: Eitan Adam Pechenick, Christopher M. Danforth, and Peter Sheridan Dodds, *Characterizing the Google Books Corpus: Strong Limits to Inferences of Socio-Cultural and Linguistic Evolution*, PLOS ONE **10** (2015), no. 10, e0137041
- Problem 1: OCR error
- The count of F-word between 1800 to 2000
- Are people suddenly become more polite after 1800?

Characteristics 9: Dirty

- “Big data sources can be loaded with junk and spam”
- Example: Eitan Adam Pechenick, Christopher M. Danforth, and Peter Sheridan Dodds, *Characterizing the Google Books Corpus: Strong Limits to Inferences of Socio-Cultural and Linguistic Evolution*, PLOS ONE **10** (2015), no. 10, e0137041
- Problem 1: OCR error
- The count of F-word between 1800 to 2000
- Are people suddenly become more polite after 1800?
- No! it's because s in old books are often written as a long s that looks like f before and around 1800s; so Google Books treat suck as the f-word in 1800.

Characteristics 9: Dirty

- Example 2: figure vs Figure

Characteristics 9: Dirty

- Example 2: figure vs Figure
- Why Figure is significantly used more than figure?

Characteristics 9: Dirty

- Example 2: figure vs Figure
- Why Figure is significantly used more than figure?
- Oversampling of scientific literature

Characteristics 10: Sensitive

- Some of the information that companies and governments have is sensitive.

Characteristics 10: Sensitive

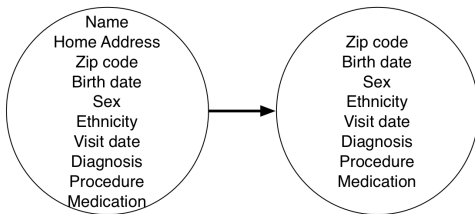
- Some of the information that companies and governments have is sensitive.
- Even if we tried to anonymize data

Characteristics 10: Sensitive

- Some of the information that companies and governments have is sensitive.
- Even if we tried to anonymize data
- This leads to potential ethical questions

Characteristics 10: Sensitive

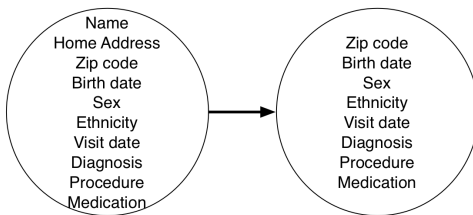
- Example 1: how anonymization fails



"Anonymization"

Characteristics 10: Sensitive

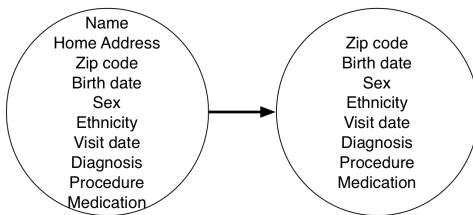
- Example 1: how anonymization fails
- Group Insurance Commission (GIC) was a government agency responsible for purchasing health insurance for all state employees in Massachusetts.



"Anonymization"

Characteristics 10: Sensitive

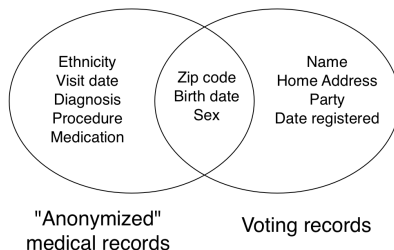
- Example 1: how anonymization fails
- Group Insurance Commission (GIC) was a government agency responsible for purchasing health insurance for all state employees in Massachusetts.
- GIC released some health information to spur research, and anonymized the part they thought were sensitive



"Anonymization"

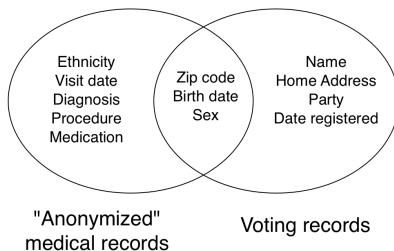
Characteristics 10: Sensitive

- But Latanya Sweeney (now a Professor at Harvard) found that she could merge GIC data with public voter registration record



Characteristics 10: Sensitive

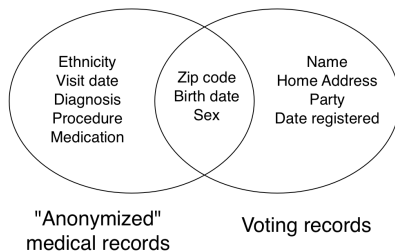
- But Latanya Sweeney (now a Professor at Harvard) found that she could merge GIC data with public voter registration record



- And in this way, she was able to find a unique match: then governor of Massachusetts.

Characteristics 10: Sensitive

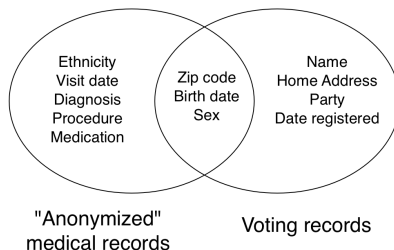
- But Latanya Sweeney (now a Professor at Harvard) found that she could merge GIC data with public voter registration record



- And in this way, she was able to find a unique match: then governor of Massachusetts.
- Sometimes, even good intention and best effort to anonymize can lead to potential harm

Characteristics 10: Sensitive

- But Latanya Sweeney (now a Professor at Harvard) found that she could merge GIC data with public voter registration record



- And in this way, she was able to find a unique match: then governor of Massachusetts.
- Sometimes, even good intention and best effort to anonymize can lead to potential harm
- Things can only be worse if no effort has been made to protect privacy

Summary

- We have summarized 10 characteristics of big data

Summary

- We have summarized 10 characteristics of big data
- You should be able to “describe the opportunities and challenges” brought by big data

Summary

- We have summarized 10 characteristics of big data
- You should be able to “describe the opportunities and challenges” brought by big data
- And when you evaluate future studies using big data

Summary

- We have summarized 10 characteristics of big data
- You should be able to “describe the opportunities and challenges” brought by big data
- And when you evaluate future studies using big data
 - **critically** evaluate their strength and weakness