

Computing in Context: Fall 2024

Lecture 1 | Overview

Second Quarter – Context-focus

- 2 lectures – more theoretical
- 1 lab – more practical
- Assignments – introduced and explained next week
- Disclaimer

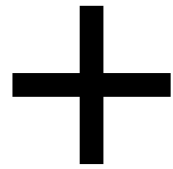
Health Policy

Health Data

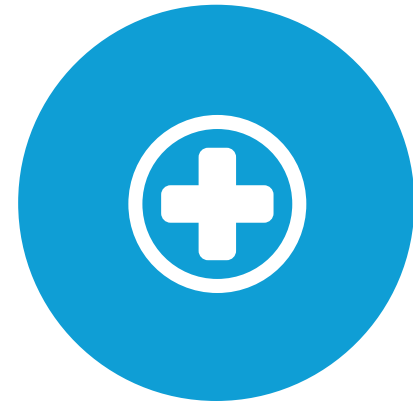
Evidence based decision making



HEALTH POLICY



HEALTH DATA

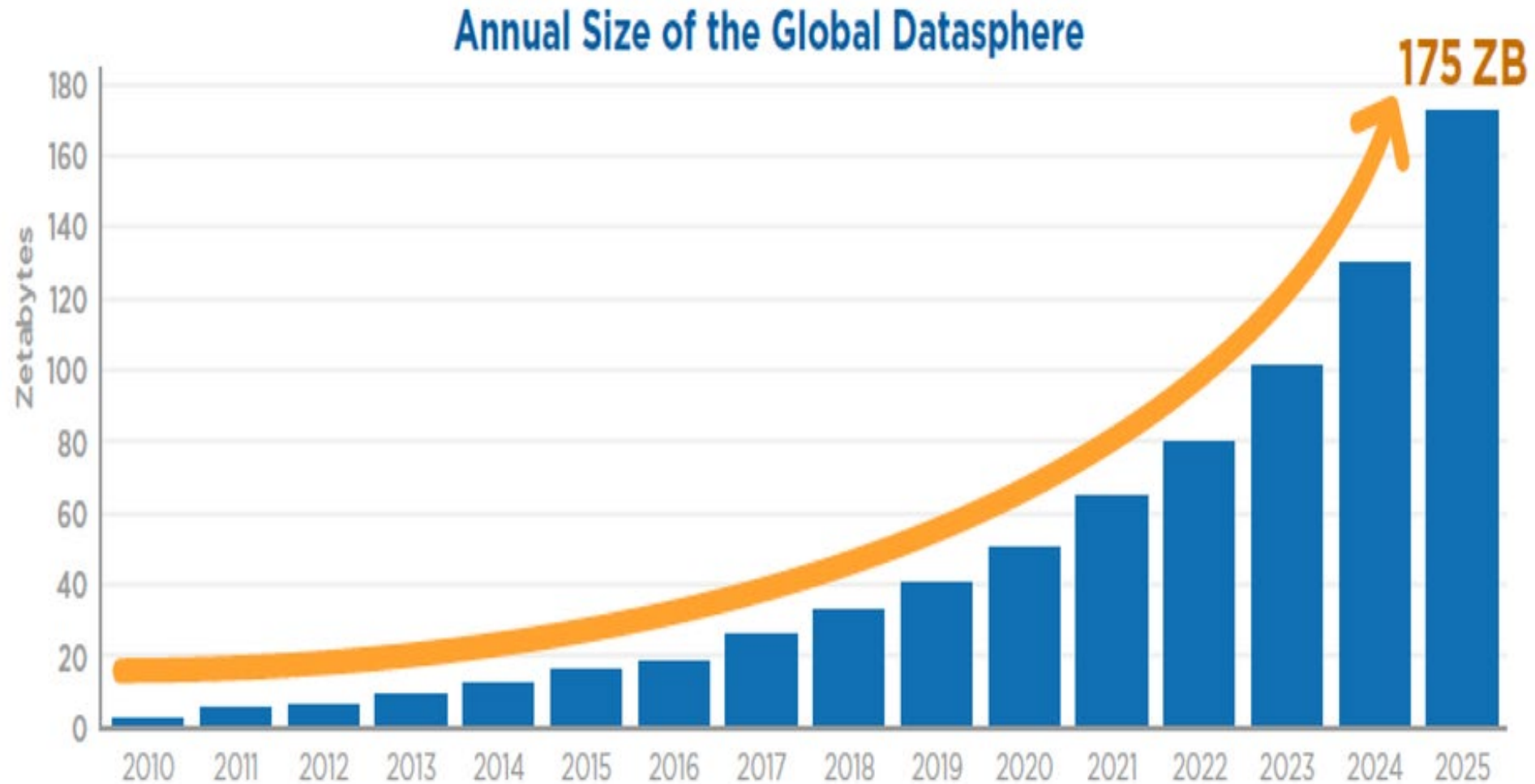


EVIDENCE BASED
HEALTH POLICY

Data and Computing

An important aside

Amount of Data Created Each Year



Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

Scale of the Global Data Growth

- 1 zettabyte = a trillion gigabytes
- 1 gigabyte = 1000 megabyte,
- 1 megabyte = 1000 kilobytes
- 1 kilobyte = 1000 bytes
- Each byte has 8 bits of information
- A bit is the smallest unit of digital information

Downloading global datasphere would take 1.8 billion years

Why does it matter?

1. Data is out there – do you need more?
2. Finding signal in the noise – how?
3. Limit may not be your data but your compute

Hardware Components of Computing

- **Processor (i.e. CPU):** The more powerful the processor is, the faster it can perform tasks assigned to it
- **Memory (i.e. RAM):** Akin to your brain's "working memory" – larger memory means the processor can operate on more info
- **Storage (i.e. HDD or SSD):** Akin to your brain's "long-term memory" – data that is not currently being used, but is kept for retrieval/editing

“The Big Data Revolution”

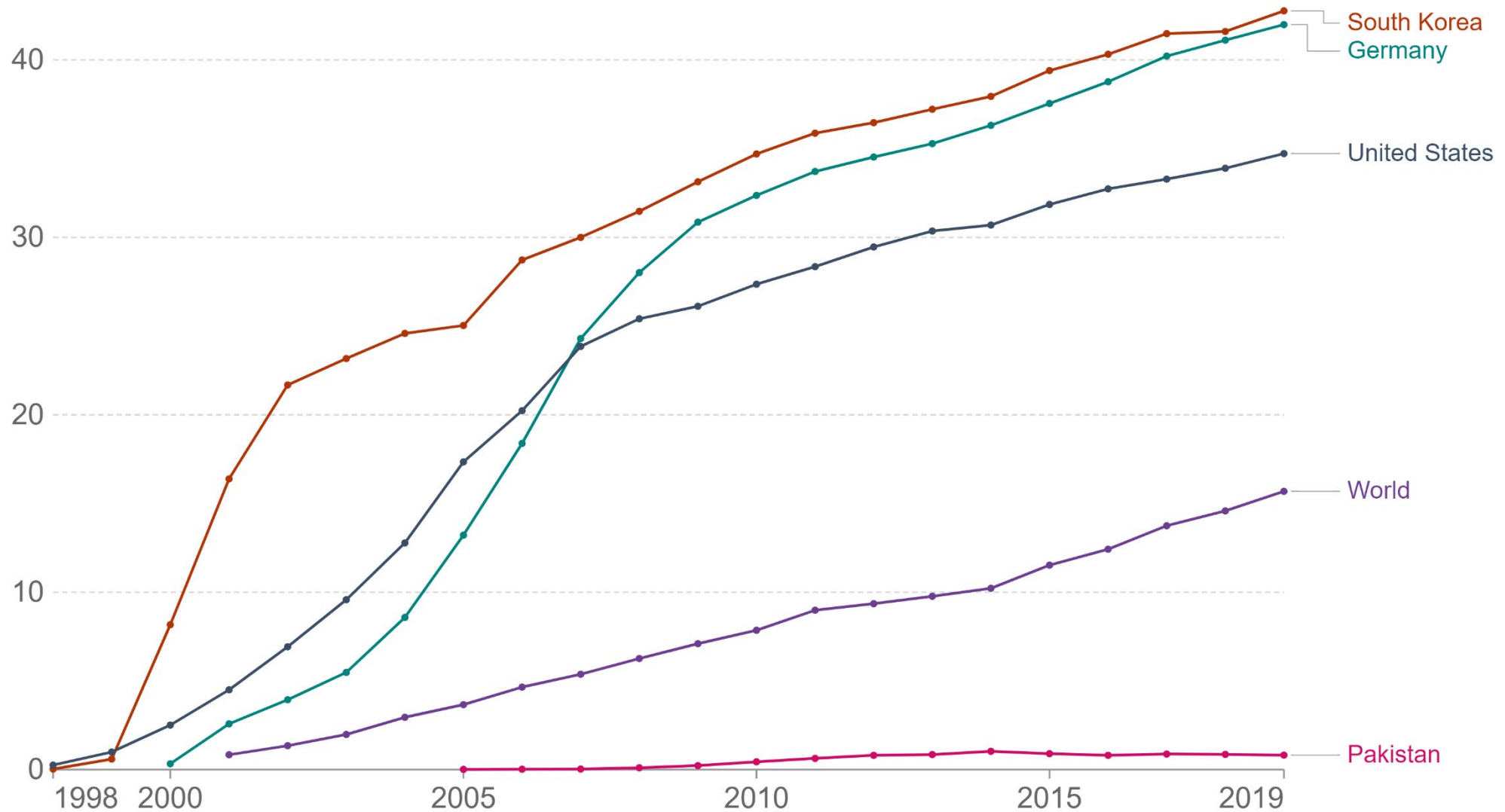
In the past 50 years - explosion in the processor speed, memory size, and storage density

In the past 20 years - broadband (high-speed) internet networks have grown, rapid transmission of large quantities of data *between computers* now feasible

- This makes data more *useful*, as it can be communicated to others

Broadband subscriptions per 100 people, 1998 to 2019

Broadband subscriptions refer to fixed subscriptions to high-speed access to the public Internet (a TCP/IP connection), at downstream speeds equal to, or greater than, 256 kbit/s.



Source: International Telecommunication Union (via World Bank)
Note: For more details on the definition see the sources tab.

Health data

Where does it come from?

Traditional Data Source: Surveys

Historically, health data not readily available (or recorded)

- Records weren't digitized (in many cases they still aren't!)

So, quantitative HP research data was collected explicitly:

- Example: surveys administered by government agencies
- e.g. the National Health Interview Survey (NHIS), administered by the Census Bureau since 1957

Upsides to survey data:

- Survey designer can choose what information is collected
- Data is generally well-structured → “ready” to analyze

Downsides to survey data:

- Expensive to collect → limits sample size & question #
- Respondents may provide inaccurate information
- Can only know what the respondent can readily recall

Administrative Data

Information on/for the operations of large organizations:

- Since Mesopotamians in 7500 BC – clay tablets for bookkeeping
- E.g. Land deeds, government proceedings, tax rolls, records of births/deaths/marriages.

From hard copy data: unreliable, errors, data loss, sharing ...



To digitization in 20th and 21st centuries with growing compute ...

Modern Administrative Data

- Accounting records
- Inventory tracking
- Customer service
- Human resources
- Tax records
- Health records
- Insurance claims

Goal of administrative data is to keep detailed records for future reference, so this data is inherently “big”:

- Both in the number of entries recorded, and the activities recorded

(Dis)Advantages of Administrative Data

Typical **advantages** of admin data relative to survey data:

- Larger sample sizes
- More precise measurement
- Less bias in responses
- Larger number of variables

Typical **disadvantages** of admin data relative to survey data:

- Not be crafted for analyst's intended purpose
- Unlikely to capture important measures of interest

Administrative Data in Healthcare

The healthcare system produces two unique, massive forms of admin data used by researchers, managers, and policymakers:

- **Health Insurance Claims**
- **Electronic Health Records (EHR)**

Sample Health Insurance Claim Form

ACME Insurance Company

HEALTH INSURANCE CLAIM FORM
APPROVED BY NATIONAL UNIFORM CLAIM COMMITTEE (NUCC) 08/12

1. MEDICARE ☒ MEDICAID ☐ TRICARE ☐ CHAMPVA ☐ GROUP HEALTH PLAN ☐ FECA ☐ OTHER ☐ 1a. INSURED'S I.D. NUMBER (For Program in Item 1)

2. PATIENT'S ADDRESS (No. Street) **Ace, Sample** 05/23/11 M ☒ F ☐ 7. INSURED'S ADDRESS (No. Street) **Ace, Sample**
111 1st street Set ☒ Spouse ☐ Child ☐ Other ☐ 111 1st street
CITY **san benito** STATE **tx** CITY **san benito** STATE **tx**
ZIP CODE **78586** TELEPHONE (Include Area Code) **(956) 7894561** ZIP CODE **78586** TELEPHONE (Include Area Code) **(956) 7894561**

3. OTHER INSURED'S NAME & ADDRESS (No. Street, City, State, ZIP Code) 10. EMPLOYMENT (Current or Previous) 11. INSURED'S DATE OF BIRTH (MM/DD/YY) SEX
a. OTHER INSURED'S POLICY OR GROUP NUMBER b. AUTO ACCIDENT? c. RESERVATION FOR NUCC USE d. INSURANCE PLAN NAME OR PROGRAM NAME 10a. CLAIM CODES (Designated by NUCC) e. IS THERE ANOTHER HEALTH BENEFIT PLAN? YES ☐ NO ☒ f. INSURED'S OR AUTHORIZED PERSON'S SIGNATURE I authorize payment of medical benefits to the undersigned physician or supplier for services described below.

12. PATIENT'S OR AUTHORIZED PERSON'S SIGNATURE I authorize the release of any medical or other information necessary to process this claim. I also request payment of government benefits either to myself or to the party who accepts assignment below. 13. INSURED'S OR AUTHORIZED PERSON'S SIGNATURE I authorize payment of medical benefits to the undersigned physician or supplier for services described below.

SIGNATURE ON FILE DATE 5/23/2011 SIGNATURE ON FILE

14. DATE OF CURRENT ILLNESS, INJURY, OR PREGNANCY (EMP) 15. OTHER DATE 16. DATES PATIENT UNABLE TO WORK IN CURRENT OCCUPATION (FROM TO)

17. NPI 34552 18. OUTSIDE LAB? YES ☐ NO ☒ 19. SUBMISSION CODE ORIGINAL REF. NO. 20. PRIOR AUTHORIZATION NUMBER

21. DIAGNOSIS OR NATURE OF ILLNESS OR INJURY (Relate A-L to service line below (24E)) ICD-9 9
A. 290.0 B. 478.6 C. 711.2 D. E. F. G. H. I. J. K. L.

22. PROCEDURE CODES (Relate A-L to service line below (24E)) CPT/HCPCS I MODIFIER POINTER \$ CHARGES UNITS \$ PAID \$ BALANCE DUE
04 11 14 04 11 14 11 92507 ABC 100.00 1 NPI 1234567890
04 11 14 04 11 14 11 92507 ABC 100.00 1 NPI 1234567890
04 11 14 04 11 14 11 92526 ABC 50.00 1 NPI 1234567890
04 11 14 04 11 14 11 92507 ABC 100.00 1 NPI 1234567890
04 11 14 04 11 14 03 92507 ABC 100.00 1 NPI 1234567890
04 11 14 04 11 14 1 92507 ABC 100.00 1 NPI 1234567890

23. FEDERAL TAX I.D. NUMBER SSN EIN 24. PATIENT'S ACCOUNT NO. 25. ACCEPT ASSIGNMENT? YES ☐ NO ☒ 26. TOTAL CHARGE 27. AMOUNT PAID 28. REVENUE FOR NUCC USE
7200000 1148 550.00 0.00
29. SIGNATURE OF PHYSICIAN OR SUPPLIER (Including Degrees or Credentials) 30. SERVICE FACILITY LOCATION INFORMATION 31. BILLING PROVIDER INFO & PH # 32. REV. FOR NUCC USE
Joe Williams, CCC SLP 123 Main Street Raleigh, NC 27609 123 Main Street Raleigh, NC 27609 1234567 22 235200000X 1234567 22 235200000X

SIGNED DATE 4/18/2014 234567 22 235200000X 1234567 22 235200000X

NUCC Instruction Manual available at: www.nucc.org PLEASE PRINT OR TYPE APPROVED OMB-0938-1197 FORM CMS 1500 (02-12)

Patient Identifiers, DOB, Sex

Provider ID, Diagnosis Codes

Procedure Codes, Date, Provider IDs

Administrative Data in Healthcare

An individual claim can't tell you much about the patient

- Basic demographics
- Diagnosis codes given by the provider
- Provider IDs, and the services they rendered

However, observing all claims for an individual over an extended period can give the analyst (you!) sense as to their health status

Electronic Health Records

On the other hand, electronic health records (EHR) can give a much clearer picture of an individual's health

- Detailed date/time data
- Patient vitals (heart rate, blood pressure, BMI, etc.), lab values
- Physician notes, medical history
- Diagnoses, Procedures, imaging (MRI, X-ray, etc.)
- And more...

Electronic Health Records (EHR)

Hyperspace - FAMSDBON PC F - WITS - KTAZD1910 WITSSDM

Desktop Action Patient Care Scheduling Billing CRM/CM Reports Report Mgmt Tools Admin Help

Back Forward Home Schedule In Basket Chart Encounter Tel Enc Message Enc Secure Panel Mgmt Print Log Out

Epic Home

Smith, John W MRN: 000017701887 Age: 30 year Sex: M PCP: Spero, Robert David (M) Allergies: Sulfia Class, Acarbose, 5-alpha Reductas Alert: Spec Fest: N kp.org: Inactive

4/22/2009 visit with TEST DUMMY MD

Images Questionnaires Admin Benefits Inquiry References SmartSets Open Orders Preview AVS Print AVS

Allergies: Sulfia Class, Acarbose, 5-alpha Reductase Inhibitors, Acetaminophen + Propoxyphene Napsylate Reviewed on 2/27/2009

Last Vitals: BP: 120/80 P: 60 T: T Src: Resp: 22 W: 190 lbs (86.183 kg) H: 5' 10" (1.778 m)
BMI: 27.26 kg/m2, BSA: 2.06 m2, Exercise Vitals: 180 mins/wk

Height: 5' 10" (1.778 m)
Peak Flow

Charting

- Chief Complaint
- Nursing Notes
- Vitals
- Exercise Vitals**
- Review Exercise VS
- Med. Document
- BestPractice
- History
- Progress Note
- SmartSets
- Dx and Orders
- Pt. Instructions
- LOS
- Follow-up
- Close Encounter

Exercise Vitals - Exercise Vitals (SHIFT+F6 to enter comments)

Instant Taken:
Date: 4/30/2009
Time: 1149

Exercise Level of Effort

Days per week of moderate to strenuous exercise (like a brisk walk): 0 1 2 3 4 5 6 7 8

On average, minutes per day of exercise at this level: 10 20 30 40 50 60 90 120 150 or greater

Restore Close F9 Cancel Previous F7 Next F8

Review Exercise Vitals

✓ Mark as Reviewed Last Reviewed by SHARMA, PANKAJ on 4/24/2009 at 12:36:26 PM

Medication Documentation

Current Prescriptions	Taking?	Start Date	End Date
ATENOLOL 100 MG ORAL TAB TAKE 1 TABLET ORALLY DAILY		4/29/2009	
ATENOLOL 100 MG ORAL TAB 1 TAB PO DAILY		4/29/2009	5/29/2011

Provider: William Lewis (M.D.) Sperling

Exit Workspace Navigator Hotkeys

Data Details

Structure, content, and “big” data

UNSTRUCTURED DATA



VS

STRUCTURED DATA



Data: “Structured” vs “Unstructured”

In your courses, it’s possible you’ve only worked with “structured” data – clearly defined and categorized easily





But most real world data is unstructured.

Intuition for “structured” vs. “unstructured”:

“Could this type of data be readily analyzed in a spreadsheet?”

“yes” = structured.

Example: Structured Data

	 person_id Scrambled individual identifier	 sample_ed Individual residing in a zip code included in the ED study	 any_visit_pre_ed Any ED visit, pre-randomization	 any_visit_ed Any ED visit in the study period
1	60562	1	0	0
2	51142	1	0	1
3	60314	1	0	1
4	50902	1	0	0
5	70733	1	0	0
6	56758	1	0	0
7	52926	1	1	1
8	4692	1	0	1
9	24115	1	0	1
10	63106	1	0	0
11	49727	1	0	0
12	44409	1	0	0
13	16262	1	0	0
14	60563	1	0	1
15	19345	1	1	0
16	69404	1	1	1
17	52156	1	1	1
18	1614	1	0	1

Healthcare Data: Structured and Unstructured

“Less than 15 percent of health data in EHRs are entered in structured data fields.” (Roski et al., 2014)

Traditionally, structure was a precondition for analysis and retrieval

But “big data” approaches enable the efficient linkage and analysis of unstructured data to answer operational or research questions

Transforming Data

A large part of what folks call “data science” consists of “ETL” (**Extract, Transform, Load**) tasks

- “How do I extract data of interest in its current (often raw, unstructured) form, and transform it so that it can be more easily used?”

Example: suppose you’d like to use the unstructured text from physicians’ EHR notes to better predict future adverse health outcomes for patients. *How might you do this?*

Example: Structuring Physician Notes

Perhaps you know a keyword that would be helpful to identify:

- Generate flag (i.e., “dummy”) variables for each patient encounter that indicate whether the physician used that word in her notes.

Maybe not and they all could be important?

- Generate a flag for each word/phrase in the physician notes?
- ... thousands of variables (think: Pandas Dataframe columns)!

Data Dimensions: “Long” vs. “Wide”

cardia c	pain	stroke
1	0	0
0	1	0
0	1	0
0	0	0
1	0	0
1	0	0
0	1	1
0	0	0
0	1	0
0	1	0
0	0	0
1	1	0
0	0	0

	cardia c	pain	stroke	male	leg	head	sharp	fever	blood
N	1	0	0	0	0	0	0	0	1
	0	1	0	1	1	0	1	0	0
	0	1	0	1	0	0	1	0	1
	0	0	0	0	0	0	0	0	1
	1	0	0	0	0	0	0	1	0
	1	0	0	0	1	0	0	0	1
	0	1	1	1	0	1	1	0	0

N: # of observations in dataset (rows)

K: # of variables in dataset (columns)

- When K is large: “wide” data
- When N is large: “long” data
- “wide” and “long” are **NOT** mutually exclusive terms!

Defining “Big Data”

“Big Data” can be used to refer to data that is wide or long, but is typically both, and often contains unstructured data fields.

Generally, we’re talking about data that’s so “big” that it cannot be contained in the memory of a typical computer.

However, most “Big Data” analysis tools (like machine learning algorithms) are uniquely suited to handling **wide** data – finding patterns across large combinations of variables.

Big Data: Linking Across Sources

Another important aspect that contributes to the ‘width’ of Big Data is that datasets from separate administrative sources are often “merged” together to address questions.

Examples:

- A health system linking SSA death records to its EHR, enabling it to better measure patient (past and present) outcomes
- Merging IRS income tax records with data from a randomized policy evaluation, to determine its effect on subsequent employment

Real-Time Digitization of Health

The advances in computing capacity have lead to proliferation of interconnected mobile devices (the “internet of things”)

This is not restricted to telecommunication and consumer electronics – consider recent developments in health

- Mobile imaging
- “Wearable” technology (heart rate and blood pressure sensors, glucose monitors, exercise trackers, health apps)

This data is all stored somewhere...

Data Misuse Concerns

Given the sensitive nature of health data, there are concerns regarding the security of personally identifiable information (PII).

It can be surprisingly easy to trace even “de-identified / anonymized” data back to a specific individual especially with the data is “wide” or the sample is specific

In addition to privacy, also concern that such data could be misused to illegally discriminate against individuals/groups

Questions?