## Plan for the week

- M: Introduction to Spreadsheets
- W: Descriptive statistics
  - Measures of central tendency
- F: Section
  - More advanced spreadsheet functionality (sort, filter, pivot tables, etc.)

# Databases vs. Spreadsheets

## Raw data can be unmanageable

Data Source

- What company appears most often in the data set?
- What nationality is most common?
- Who is the youngest person in the dataset?
- What is the average net worth of all the people in the dataset?

name	networth	nationality	company	age
Carlos Slim Helu	\$54.5 Billion	Mexico	América Móvil	77
David Koch	\$48.3 Billion	U.S.	Koch Industries	77
Charles Koch	\$48.3 Billion	U.S.	Koch Industries	81
Bill Gates	\$86 Billion	U.S.	Microsoft	60
Sheldon Adelson	\$30.4 Billion	U.S.	Las Vegas Sands	83
Warren Buffett	\$75.6 Billion	U.S.	Berkshire Hathaway	85
Sergey Brin	\$39.8 Billion	U.S.	Google	43
Larry Ellison	\$52.5 Billion	U.S.	Oracle	72
Wang Jianlin	\$31.3 Billion	China	Dallan Wanda Group	62
Li Ka-shing	\$31.2 Billion	Hong Kong	CK Hutchison Holding	89
Liliane Bettencourt	\$39.5 Billion	French	L'Oreal	94
Larry Page	\$39.8 Billion	U.S.	Google	43
Amancio Ortega	\$71.3 Billion	Spain	Zara	81
Jim Walton	\$34 Billion	U.S.	Walmart	60
Rob Walton	\$34.1 Billion	U.S.	Walmart	72
Bernard Arnault	\$41.5 Billion	French	Louis Vuitton	68
Alice Walton	\$33.8 Billion	U.S.	Walmart	67
Jeff Bezos	\$72.8 Billion	U.S.	Amazon	53
Mark Zuckerberg	\$56 Billion	U.S.	Facebook	33
Michael Bloomberg	\$47.5 Billion	U.S.	Bloomberg	75

## Raw data can be unmanageable

- What company appears most often in the data set?
  - Walmart
- What nationality is most common?
  - **U.S.**
- Who is the youngest person?
  Mark Zuckerberg
- What is the average networth of the people in the dataset?
  - \$48.41 Billion

name	networth	nationality	company	age
Carlos Slim Helu	\$54.5 Billion	Mexico	América Móvil	77
David Koch	\$48.3 Billion	U.S.	Koch Industries	77
Charles Koch	\$48.3 Billion	U.S.	Koch Industries	81
Bill Gates	\$86 Billion	U.S.	Microsoft	60
Sheldon Adelson	\$30.4 Billion	U.S.	Las Vegas Sands	83
Warren Buffett	\$75.6 Billion	U.S.	Berkshire Hathaway	85
Sergey Brin	\$39.8 Billion	U.S.	Google	43
Larry Ellison	\$52.5 Billion	U.S.	Oracle	72
Wang Jianlin	\$31.3 Billion	China	Dallan Wanda Group	62
Li Ka-shing	\$31.2 Billion	Hong Kong	CK Hutchison Holding	89
Liliane Bettencourt	\$39.5 Billion	French	L'Oreal	94
Larry Page	\$39.8 Billion	U.S.	Google	43
Amancio Ortega	\$71.3 Billion	Spain	Zara	81
Jim Walton	\$34 Billion	U.S.	Walmart	60
Rob Walton	\$34.1 Billion	U.S.	Walmart	72
Bernard Arnault	\$41.5 Billion	French	Louis Vuitton	68
Alice Walton	\$33.8 Billion	U.S.	Walmart	67
Jeff Bezos	\$72.8 Billion	U.S.	Amazon	53
Mark Zuckerberg	\$56 Billion	U.S.	Facebook	33
Michael Bloomberg	\$47.5 Billion	U.S.	Bloomberg	75

Data Source

## Raw data can be unmanageable

Data Source

- It is often impossible to answer these or similar queries, or to establish any trends, by eyeballing the data.
- Computational processing is necessary for data beyond a trivial size.

name	networth	nationality	company	age
Carlos Slim Helu	\$54.5 Billion	Mexico	América Móvil	77
David Koch	\$48.3 Billion	U.S.	Koch Industries	77
Charles Koch	\$48.3 Billion	U.S.	Koch Industries	81
Bill Gates	\$86 Billion	U.S.	Microsoft	60
Sheldon Adelson	\$30.4 Billion	U.S.	Las Vegas Sands	83
Warren Buffett	\$75.6 Billion	U.S.	Berkshire Hathaway	85
Sergey Brin	\$39.8 Billion	U.S.	Google	43
Larry Ellison	\$52.5 Billion	U.S.	Oracle	72
Wang Jianlin	\$31.3 Billion	China	Dallan Wanda Group	62
Li Ka-shing	\$31.2 Billion	Hong Kong	CK Hutchison Holding	89
Liliane Bettencourt	\$39.5 Billion	French	L'Oreal	94
Larry Page	\$39.8 Billion	U.S.	Google	43
Amancio Ortega	\$71.3 Billion	Spain	Zara	81
Jim Walton	\$34 Billion	U.S.	Walmart	60
Rob Walton	\$34.1 Billion	U.S.	Walmart	72
Bernard Arnault	\$41.5 Billion	French	Louis Vuitton	68
Alice Walton	\$33.8 Billion	U.S.	Walmart	67
Jeff Bezos	\$72.8 Billion	U.S.	Amazon	53
Mark Zuckerberg	\$56 Billion	U.S.	Facebook	33
Michael Bloomberg	\$47.5 Billion	U.S.	Bloomberg	75

### Databases

- A database is a structured set of data that are easily accessible in various ways
- Database software tools facilitate the automated management of data
  - Storing, searching, modifying, and extracting information in a database
- Database systems can manage a very large quantity of data

## Database Software

- Database Management Systems: software that handles storage, retrieval, and updating of databases
  - Examples: Oracle, MySQL, Access
- Statistics Packages
  - Examples: SAS, SPSS, Stata, R
- Spreadsheets: they do a little of both, at human scale, and for human comprehension, with minimal human effort!

## Spreadsheets

- Suitable for managing relatively small databases
- Designed to facilitate human comprehension of data
- Useful for visualizing both raw data, in tabular form, and summaries of data in the form of charts, etc.
- Data management is usually accomplished manually, not via automated tools

# History of Databases

## Public School Database then

- Before computers were commonplace, people used other systems to track student data
- Schools often used paper-based books (can you imagine?!) for grades, attendance, student contact info, etc.
- These books would be divided into many small squares where teachers would record relevant information

							-		1	91	18	
				(	Oc	to	be	2				
	MT	Wq	512	M 14	W' 16	5	M 21	23.	526	28	W 30	
& R. Stearn	V	V	1	+	79	~	hach	tod	. 1	tes	A.B.	
7. Lanaridae	V	V	1	V.	V	V	V	V	V	V	V	
C.W. Tompinson	V	V	1	ď	d	4	47			-		1
S.W. Jurtle"	V	V	×	L	X	x	d	V	V	V	V	0
C. Schiff	V	V	V	V	V	V	V	V	×	d	d	4
P. Steinberg	V	1	V	V	V	V	V	d	d	d	d	0
H.S.N. menko	V	Ľ	L	V	V	V	d	d	x	x	X	(
A. Gullestein	V	1	V	V	V	V	V	1	V	1	1	4
R.N. Downing	V	V	d	V	V	V	d	d	V	d	1	1
D.R. Thompson	v	V	d	d	d	d	L	4	L	d	d	-

## Difficulties with paper-based records

- Mistakes are hard to correct and quickly become quite messy
- It is difficult to search through the whole book (slowly or quickly)
- Calculating summary statistics is difficult (e.g., how many days did R.N. Downing miss school?)
- Records can be lost forever or damaged due to water, fire, etc., because it is expensive to store backups

1- Charlester									1	91	18	
				(	Oc	to	be	2				
	M	Wq	512	M 14	W 16	5	M 21	23	526	m. 28	W 30	
E. R. Stearn	V	V	1	+	79	~	hach	tod	10 %	tes	.B.	
J. Langridge	V	V	1	V.	V	V	V	V	V	V	V	1
C.W. Tompinson	V	V	1	ď	d	4	t t			-		
S.W. Jurtle"	V	V	×	L	d	x	d	V	V	V	V	0
C. Schiff	V	V	V	V	V	V	V	V	*	d	d	c
P. Steinberg	V	1	V	V	V	V	V	d	æ	d	d	0
H.S.N. menko	V	Ľ	L	V	V	V	ď	d	x	or	X	(
A. Gullestein	V	1	V	V	V	V	V	1	V	1	1	0
R.N. Downing	V	V	d	V	V	V	d,	d	1	d	1	-
D.R. Thompson	v	V	d	d	d	ď	L	4	L	d	d	-

Image Source

## Public School Database now

- Make heavy use of (hopefully, secure) spreadsheets
- Mistakes are easy to fix (and make!) in spreadsheets
  - One can simply overwrite an entry
- Can easily search for specific entries
  - E.g., did R.N. Downing attend school on 10/14?
- Spreadsheets can also be used to easily calculate summary statistics, like grade-point averages and total number of days absent
- Records are backed up on some online system, so they are less likely to be damaged like paper books

1000	Α	В	С		D	E	F	G
7	Stu	dent Name	Monday	T	uesday	/ Wednesday	Thursday	Friday
8	314		15 Nov-15	16	6 Nov-15	5 17 Nov-15	18 Nov-15	19 Nov-15
9	01	Lee, John	PR		PR	PR	PR	PR
10	02	Williams, Adam	PR		PR	PR	PR	PR
11	03	Williams, Sarah	PR		PR	PR	PR	PR
12	04	Doe, Jane	PR	8	AB	AB	PR	AB
13	05	Liu, Kat	PR		AB	PR	PR	PR
14	06	Smith, Elizabeth	PR	8 	PR	PR	PR	PR
15	07	Sanders, Shane	AB		PR	PR	PR	PR
16	08	White, Andres	PR	8	PR	PR	PR	PR
17	09	Aaron son, Allis on	PR		PR	PR	PR	PR
18	10	Mapple, Noah	PR	8 	PR	PR	PR	PR
19	11	Grande, John	PR		PR	PR	PR	PR
20	12	McDonell, Josh	PR	8	PR	PR	PR	PR
21	13	Spear, Terry	PR		PR	PR	PR	PR
22	14	Greene, Tom	PR	8	PR	PR	PR	PR
23	15	Snow, Daniel	PR		PR	PR	PR	PR
24	16	Daniels, Jack	PR	8	PR	AB	PR	PR
25	17	Daniels, Bob	PR	0.08 0.08	PR	PR	PR	PR
26	18	West, Williams	PR	±	PR	PR	PR	PR
27	19	10 B			14			
28	20			8	A	TENDANCE:		
29	21				S	elect P if Present		
30	22			8	A	B if Absent		
31	23							
32	24			8				
33	25							
34	26			8				
35	27							
36	28			8				
37	29							
38	30			8				
39	H	N° PResent per Day:	17		16	16	18	17
40		Nº ABsent per Day:	1	8	2	2	0	1
41	18	Total:	18		18	18	18	18
42	11	Daily Attendance %:	94	8	89	89	100	94

	A	В	С	D	E	F	G	Н	1	
1				() ()		Weekly Atte	ndance Recor	d		8 - S
2										
3	Tea	cher:		Ms. S	Smith	ar 3	Group Name	4 <sup>th</sup> Grade		· - 33 -
4	WE	EK BEGINNING: (dd/mm/gggg)		11/15	/2015		Level :			
5	· · · · ·			Sun	nday					
6			0.000	i n s i	10000 A.	i na sa i			Longer	il.
7	Stu	dent Name	Monday	Tuesday	Wednesday	Thursday	Friday	Notes	ATTEN	DANCE
8	-	las lab	15 NOV-15	16 Nov-15	17 NOV-15	18 Nov-15	19 Nov-15		P Resent	Attent
10	01	Lee, John	PR	PR	PR.	PR	PR	Wail in hour sale on Weicheelday (Liveral)	5	0
10	02	Williams, Adam	PR	PR	PR	PR	PR		5	0
10	03	volitans, salat	PR	PR	PR	PR	PR		5	0
12	04	Doe, Jale	PR.	AD	AD	PR	AD		2	
13	05	LII, Fat	PR	AB	PR	PR	PR		+	
14	06	Smith, Elizabeth	PR	PR	PR	PR	PR		5	0
15	07	Salders, Shake	AB	PR	PR	PR	PR		+	
10	08	White, Andres	PR	PR	PR	PR	PR		5	0
1/	09	Aaroisol, Alisol	PR	PR	PR	PR	PR		5	0
18	10	Mappe, Noa	PR	PR	PR	PR	PR		5	0
19	11	Graide, John	PR	PR	PR	PR	PR		5	0
20	12	MCDOLEII, JOS L	PR	PR	PR.	PR	PR		5	0
21	13	Spear, Terry	PR	PR	PR	PR	PR		5	0
22	14	G ree i e , Tom	PR	PR	PR.	PR	PR		5	0
23	15	Stow, Datiel	PR	PR	PR.	PR	PR		5	0
24	16	Daile is, Jack	PR	PR	AB	ATTE	NDANCE:		+	1
25	17	Davie is, Bob	PR	PR	PR.	Sele	rt ·		5	0
26	18	West, Williams	PR	PR	PR.	PR if	Present		5	0
27	19					AB if	Absent		0	0
28	20			i i			Abaciic		0	0
29	21								0	0
30	22								0	0
31	23			1					0	0
32	24								0	0
33	25			8		6			0	0
34	26			a					0	0
35	27			1					0	0
36	28			1					0	0
37	29								0	0
38	30			8		33			0	O
39		N* <u>PResent</u> per Day:	17	16	16	18	17	Week Attendance Totals:	84	8
40	18	N* <u>ABcent</u> per Day:	1	2	2	0	1			
41	, "I	To tal :	18	18	18	18	18	Days Recorded:	6	6 6
42		Dally Attendance %:	94	89	89	100	94	Week Atlendance %:	93	7

## Summarizing the Data

Number of students in attendance per day

	A	В	C	D	E
1	Data		53 		
2	Sum · Mond≉	Sum ∙Tuesd≯	Sum - Wedne	Sum - Thurso	Sum · Friday
3	17	18	16	16	17

- Student with the most absences: Jane Doe
- Student with the most Monday absences: Shane Sanders
- Etc.

## (Coarse) Visualizations



#### Number of students in attendance per day

Percentage of students in attendance per day



History of Spreadsheets

#### History of Spreadsheets

- The word spreadsheet comes from the word "spread," as in a newspaper or magazine
- "Spreadsheet" was used to refer to bookkeeping ledgers with all revenue, costs, taxes, etc. spread on a single sheet of paper (or across two pages of a bound ledger)
- In 1961, Professor Richard Mattessich pioneered computerized spreadsheets for business accounting on main frames: "Spreadsheet: Its First Computerization (1961-1964)"

#### VisiCalc: Dan Bricklin & Bob Frankston

Subject: RE: Spreadsheets history Date: Wed, 12 May 1999 17:25:22 -0400 From: Dan Bricklin Organization: Trellix Corporation To: 'Dan Power' CC: "'Bob\_Frankston@frankston.com'"

Dan,

As I said I'd do in my email a few weeks ago (I've included a copy below) I've posted on my web site a first attempt at addressing what's so special about VisiCalc that makes people call it the "first" electronic spreadsheet as we know them today. As to the story of coming up with the idea and refining it into a product, I have the beginnings on my site (which you've read already, I assume) and Bob and I are working on more information about design decisions, exactly what happened, etc. Since we were often the only ones there, I hope you'll put some trust in our story (which we'll try to back up with pictures, check with friends from that time, etc.). I'll be seeing some of my classmates from Harvard in a few weeks and will try to get more first-hand accounts of the history.

Here's the URL for my posting:

If you want some nit picks on your piece (if not, skip this paragraph): I programmed the first "working prototypes" (not "working version") in the fall of 1978 (Oct? Nov? I'll have to check -- I think it was over a long weekend), not the summer. I recruited Bob to do a real, assembler version, instead of the prototype in Integer Basic. The prototypes actually had a better interface in some cases (context sensitive help, for example). Bob was to build production code (faster speed, better arithmetic, scrolling, etc.). Since one of the things that made VisiCalc special was the implementation and the details only worked out in that implementation (some of which we will cover in that material Bob and I are working on) Bob is referred to as the "co-creator" of VisiCalc. If you want to include his many contributions in the category of invention, then he certainly is a "co-inventor", but there is no question he was a co-creator. Fylstra, Bob, and I first talked about VisiCalc in the fall (after summer vacation I call fall) of 1978. I'll get the actual dates of some of these meetings later (for the record, not that many people care...). I don't think the number of sold copies ever made it over several hundred thousand, but I don't have the records ("about 1 million" sounds nice and if you round 500K+ up, then it's probably correct). [Enough nits... Tell me if you want things at this detail. :)]

Let me know what you think.

Thanks again for your interest in the history of this area.

#### VisiCalc

# Big idea: Instant automatic recalculation based on formulas stored in the cells referencing other cells



**Image Source** 

#### Lotus 1-2-3

- January 1983
- Mitch Kapor and Jonathan Sachs
- More powerful than VisiCalc
- Very popular
- Lotus 1-2-3 is thought to be one of the reasons the IBM PC was so successful. Likewise, for VisiCalc and the Apple II.

A : A:	1: 'EMP		97 - 1970 - 19	M 2 21 10		101 ale	5	MENU
Wor	ksheet R	ange Copy	Move File	Print G	raph Data	System	Quit	
Glo	bal Inser	rt Delete	Column Eras	e Title:	s Window	Status	Page Hide	
Ĥ	A	B	C	0	E	F	G	
1	EMP	emp_name	DEPTNO	JOB	YEARS	SALARY	BONUS	
2	1777	Azibad	4000	Sales	2	40000	10000	
3	81964	Brown	6000	Sales	3	45000	10000	
4	40370	Burns	6000	Mgr	4	75000	25000	
5	50706	Caeser	7000	Mgr	3	65000	25000	
6	49692	Curly	3000	Mgr	5	65000	20000	
7	34791	Dabarrett	7000	Sales	2	45000	10000	
8	84984	Daniels	1000	Preside	nt 8	150000	100000	
9	59937	Dempsey	3000	Sales	3	40000	10000	
10	51515	Donovan	3000	Sales	2	30000	5000	
11	48338	Fields	4000	Mgr	5	70000	25000	
12	91574	Fiklore	1000	Admin	8	35000		
13	64596	Fine	5000	Mgr	3	75000	25000	
14	13729	Green	1000	Mgr	5	90000	25000	
15	55957	Hermann	4000	Sales	4	50000	10000	
16	31619	Hodgedon	5000	Sales	2	40000	10000	
17	1773	Howard	2000	Mgr	3	80000	25000	
18	2165	Hugh	1000	Admin	5	30000		
19	23907	Johnson	1000	VP	1	100000	50000	
20	7166	Laflare	2000	Sales	2	35000	5000	
DATI	A WK3							

#### **Microsoft Excel**

- Mid 1980's
- First spreadsheet to use a GUI, or graphical user interface, meaning not an entirely text-based interface
- Was a big selling point for Apple's Macintosh; people bought it because it came with Excel
- In 1987, Microsoft released Windows with Excel

#### Today there are many options

- Microsoft Excel
- Google Spreadsheets
- LibreOffice Calc
  - The screenshots in these slides are of LibreOffice Calc, because it used to be used in CS 100

$\begin{array}{c c c c c c c c c c c c c c c c c c c $	an :↓ :1   • ● % m an an G	● <b>●</b>   ◆ 9   • • • • • • • • • • • • • • • • • • •		-   🗐	К	
Anal  10  A  A  E  III    •  #  Σ  Ξ	5 G	H	1		к	
A B C D E F	G	н	I	J	К	
	G	н	1	J	К	
						-
						-
	_					-
						-
FF Sheet1 (+ m	1					

# Qualitative vs. Quantitative Data

## Data can be either qualitative or quantitative

Qualitative Quantitative Like Easy Awkward <sub>Slow</sub> 23,406 2m32s 76.8% Squirrel Efficient 45,849 Ambiguous HOW 1,127 3.76% Confusing €12.75

Image Source

## Qualitative data

Qualitative data describe qualities, like color, texture, smell, taste, appearance, etc.

Many qualitative data are categorical: e.g.,

- the color of a ball (yellow, blue, or red)
- the brand of a product purchased (brand A, B, or C)
- whether a person is employed (yes or no)





## Qualitative data can be nominal or ordinal

• Nominal means that there is no natural order among the values





Nominal

Cows

Dogs

Pias

## Quantitative data

Quantitative data take on numerical values, so are typically ordinal

Examples:

- age, weight, height, income, etc.
- the value of a country's exports
- a batter's number of home runs



**Image Source** 

#### Quantitative data can be either discrete or continuous

• Data are discrete if the measurements are necessarily integral (i.e., integers)

 Data are continuous if the measurements can take on any value, usually within some range



Discrete

•• 5