# Data Visualization Discover, Explore and be Skeptical

Di Cook Statistics, Iowa State University soon to be Business Analytics, Monash University

# Seminar 1 Exploratory Data Analysis



Les Diablerets, Feb 1-4, 2015



Les Diablerets, Feb 1-4, 2015

DATA

Results from 2012 OECD PISA tests: "the world's global metric for quality, equity and efficiency in school education". 485,490 students math, science and reading test scores 65 countries, between 100-1500 schools in each Student questionnaires about their environment (635 vars) Parents surveyed on work, life, income (143 vars)

Principals provide information about their schools (291 vars)

### http://www.oecd.org/pisa/pisaproducts/ datavisualizationcontest.htm



-40 Reading Score Gap In reading scores, the world is PINK! In EVERY COUNTRY tested GIRLS score significantly BETTER THAN BOYS. Why don't we talk about the gender gap in reading?

Les Diablerets, Feb 1-4, 2015

"Show the data."

9 -73

"The greatest value of a picture is when it forces us to notice what we never expected to see."

Examples

- OECD PISA 2012
- CO<sub>2</sub> Emissions and Climate Change
- Sports: tennis

### Data analysis







14 -73

### EDA vs IDA

- EDA and IDA (Chatfield; Crowder & Hand), although not entirely distinct, differ in emphasis.
- Fundamental to EDA is the desire to let the data inform us, to approach the data without pre-conceived hypotheses, so that we may discover unexpected features. Of course, some of the unexpected features may be errors in the data. IDA emphasizes finding these errors by checking the quality of data prior to formal modeling.
- IDA is much more closely tied to inference than EDA: Problems with the data that violate the assumptions required for valid inference need to be discovered and fixed early.

# Data snooping

- Because EDA is very graphical, it sometimes gives rise to a suspicion that patterns in the data are being detected and reported that are not really there.
- This is called data snooping.
- Validation of findings, by all means possible, is an important component of data analysis.

### Data snooping

"In our experience, false discovery is the lesser danger when compared to nondiscovery. Nondiscovery is the failure to identify meaningful structure, and it may result in false or incomplete modeling. In a healthy scientific enterprise, the fear of nondiscovery should be at least as great as the fear of false discovery." (Buja's remark on discussion on Koschat & Swayne (1996))

### PISA 2012 Process

- Files distributed as an external representation of an R object (.rda file), dictionary files=variable coding.
- Read in all files to R, browse size of files, dictionary of items, make some tables/pictures
- Make a list of questions to explore
- Read map data from rworldmap package, coordinate names for all countries in OECD data, and merge polygons with subset of variables

17 -73

Example

#### > head(dict student2012, 20)

	•		
description	able	1	
Country code 3-character	CNT	L	
sub-region code 7-digit code (3-digit country code + region ID + stratum ID)	ATIO	2 8	
Stratum ID 7-character (cnt + region ID + original stratum ID	ATUM	3	
OECD country	OECD	1	
National Centre 6-digit Code	NC	5	
School ID 7-digit (region ID + stratum ID + 3-digit school ID	OLID	5	
Student II	DSTD	7	
International Grade	1Q01	3	
National Study Programme	2Q01	9	
Birth - Month	3Q01	L 0	
Birth -Year	3Q02	11	
Gender	4Q01	12	
Attend <isced 0=""></isced>	5Q01	13	
Age at <isced 1=""></isced>	6Q01	L4	
Repeat - <isced 1=""></isced>	7Q01	15	
Repeat - <isced 2=""></isced>	7Q02	L 6	
Repeat - <isced 3=""></isced>	7Q03	L7	
Truancy - Late for School	8Q01	L 8	
Truancy - Skip whole school day	9Q01	19	
Truancy - Skip classes within school day	5Q01	20 8	

Example

<pre>&gt; tail(dict_student2012</pre>	2, 5)
variable	description
631 W_FSTR80	FINAL STUDENT REPLICATE BRR-FAY WEIGHT80
	DANDONTZED ETNAT VADTANCE CODADIM (1 90)

	_			
632	WVARSTRR		RANDOMIZED FINAL VARIANCE STRATUM (1-80)	)
633	VAR_UNIT		RANDOMLY ASSIGNED VARIANCE UNIT	Г
634	SENWGT_STU	Senate weight -	sum of weight within the country is 1000	)
635	VER STU		Date of the database creation	ı



Compare plausible values

### Example

		lent2012\$CNT))	sort(table(stud
Perm(Russian Federation)	Massachusetts (USA)	Connecticut (USA)	Liechtenstein
1761	1723	1697	293
Latvia	New Zealand	Iceland	Florida (USA)
4306	4291	3508	1896
Poland	Costa Rica	Netherlands	Tunisia
4607	4602	4460	4407
Slovak Republic	Hong Kong-China	Lithuania	France
4678	4670	4618	4613
		•••	
Uruguay	Bulgaria	Luxembourg	Russian Federation
5315	5282	5258	5231
Indonesia	Singapore	Macao-China	Czech Republic
5622	5546	5335	5327
Slovenia	Argentina	Kazakhstan	Portugal
5911	5908	5808	5722
Thailand	Japan	Chinese Taipei	Peru
6606	6351	6046	6035
Belgium	Denmark	Jordan	Chile
8597	7481	7038	6856
Switzerland	Qatar	Colombia	Finland
11229	10966	9073	8829
Brazil	Australia	United Kingdom	United Arab Emirates
19204	14481	12659	11500
Mexico	Italy	Spain	Canada
33806	31073	25313	21544

Les Diablerets, Feb 1-4, 2015

Compare plausible values for Switzerland



Do sample weights matter?



### Questions

- Is the gender gap evident in the math scores?
- Do students work hard outside school hours? And does this pay off in better scores?
- Does truancy affect performance?
- Do more household possessions mean better performance?
- Do parents matter?

...

LES DIABLERETS, FEB 1-4, 201

25 -73

# Look at the individuals, not just summary statistics



### Expectations

- @ Gender gap: "boys are better at math"
- ISA doesn't perform well compared to other countries.
- I would expect parents, time studying and socioeconomic status to be important influences on performance.



- Test scores range from 0-1000.
- Gaps for math were at most 30 points.
- For reading at most 80 points.
- Differences are tiny.

LES DIABLERETS, FEB 1-4, 2015



# Individuals

• Even countries with a math gender gap, sometimes have a girl who scored the top mark.

### Czech Republic

29-73



Hours spent out studying out of school time and 600 500 400 700 600 500 9400 400 600 500 average math score by 400 700 country 600 500 400





# Individuals

• And although there is a gender gap in reading in all countries, individually some top reading scores are obtained by boys.





### Reported truancy and average math score by country



33-73

# Parents in the home and average math score by country





Family structure



# Possessions and average math score by country



37-73

39-73

### Number of TVs and average math score by country





### Reported internet use by Swiss teens



What calculations, tables, plots would you make to tackle these questions?

### What we learned

- @ Gender gap is not universal in math.
- @ Gender gap is universal in reading.
- Individual differences trump everything else.
- Parents matter! Mothers more than fathers!
- Socioeconomic status matters.
- Turn the TV off, if you live in a developed country!
- 🚱 Something's funny about Albania's data.

Les Diablerets, Feb 1-4, 2015

### Graphics choices

Sorting!

- Ordered dot plots
- Ine plots
- Histograms
- Showing stats vs all values
- Our Counts of Category Counts

		Climate change				
		CO2 alarm bells?	"Planet's CO <sub>2</sub> l 400 ppm for fi human exi	evel reaches rst time in stence."		
		"It's tec questi	s not a question of if air hnology will be adopted ion of when," said Klaus	capture l; it's a s Lackner		
		"In the last 80 years, as CO <sub>2</sub> em annual rate of climate-related dea This means the incidence of dear was 80 years	vissions have most rapidly e aths worldwide fell by an in th from climate is 50 times ago," writes Epstein.	scalated, the acredible 98%. lower than it		
Les Diablerets, Feb 1-4, 2015	45 -73	Les Diable	rets, Feb 1-4, 2015	46 -73		
Climate change		Clima	te change			

### Data available at <u>http://scrippsco2.ucsd.edu</u>

> head(CO2.all)											
		date	time	day	decdate	n	flg	co2	lat	lon	$\mathtt{stn}$
1	1	1997-01-21	13:40	35451.57	1997.056	4	0	365.82	23.3	-110.2	bcs
1	2	1997-02-08	15:03	35469.63	1997.106	2	0	365.57	23.3	-110.2	bcs
;	3	1997-02-22	15 <b>:</b> 07	35483.63	1997.144	2	0	366.28	23.3	-110.2	bcs
	4	1997-03-07	14:23	35496.60	1997.180	3	0	367.85	23.3	-110.2	bcs
1	5	1997-03-22	14:16	35511.59	1997.221	2	0	365.28	23.3	-110.2	bcs
	б	1997-04-05	13 <b>:</b> 37	35525.57	1997.259	3	0	368.96	23.3	-110.2	bcs

I expected:

© CO2 concentrations to be flattening over time

🚱 Raw data to look like raw data - messy







Minimum and maximum for each year and station, show that they are effectively measuring the same product



For station, ptb, most northerly station, examine monthly trend. Drop in CO<sub>2</sub> occurs every year in June.

Strange year is 1978, and lack measurements for summer.

Climate change

Surprises:

- © CO2 concentrations increasing consistently.
- Worthern stations show seasonality.
- Oata looks very regular. Every recording station essentially producing same values, ignoring seasonality, with a few slight anomalies.



Residuals from a linear model fit indicate the rate of increase is exponential.



Change aspect ratio and anomalies at different sites are visible.



### Climate change

What's happening near Les Diablerets?

- 🕼 Global Historical Climatological Network data
- 🕼 Station: Col du Grand St Bernard
- Minimum/maximum temperature and precipitation



les Diablerets, Feb 1-4, 2015

57 -73



### Maximum daily temperature

Seasonal trend visible from plotting temperature by month Temperature does NOT reach freezing in SUMMER occasionally





### Maximum daily temperature

Plot of temperature by year reveals a big gap in measurements!

May need to focus on one period or the other.





62 -73

Change aspect ratio for studying

season effects

Precipitation (sqrt scale)

3 2

Small seasonal trend visible from plotting precip by month increase Feb-Jun.

temperature?



Why are there so many missing values in the last decade?



### Precipitation (sqrt scale)

Small gap in measurement collections. Change to only look at latter measurements.

And compute monthly totals (adjusted for missing values).















### Precipitation (sqrt scale)



Are extremes becoming more common?

66 -73

# Your Turn

- Based on what we have seen with min/max temperature and precip what might you do to find additional information that supports or refutes?
- What data, calculations, plots would you make?

#### les Diablerets, Feb 1-4, 2015

LES DIABLERETS, FEB 1-4, 20

### 72 -73

Pulling data to check things for yourself is so much easier today. One of the most useful skills you can attain!

Summary

Graphics for exploratory data analysis are ephemeral. Once created, once things are learned, they evaporate, and need to be replaced with something overly produced and beautiful for general consumption.

Les Diablerets, Feb 1-4, 2015

69 -73

Acknowledgements

Code is available at:

https://github.com/dicook/lesdiablerets-code

Software used: R (<u>http://www.R-project.org</u>) and primarily the package ggplot2

### Coming next...

- Do you like tennis?
- Data doesn't come with only two variables. How do you see in high dimensions?
- It was the set of t

Les Diablerets, Feb 1-4, 2015