NZAMT 2013

Caught in the Path of the Data Deluge

Chris Wild

Department of Statistics University of Auckland, New Zealand



Laying story-telling foundations



Comparing heights of boys and girls at age 12



Aside: Marking dot plot up with with boxplot



Comparing heights of boys and girls at age 12



Learning to see shifts

E UNIVERSITY OF AUCKL

DEPARTMENT OF STATISTICS

Learning to see shifts



THE UNIVERSITY OF AUCKLAND DEPARTMENT OF STATISTICS





Very Useful

Effect of Sample Size









Caught in the Path of the Data Deluge

Chris Wild

Department of Statistics University of Auckland, New Zealand



Looking downwards

NCEA Level 8 Stats 2013

 Getting my head around teaching all this new stuff – it's a lot of work

How to play the new NCEA games well? (for my kids)

What are the new game's rules??? Don't seem clear yet *!!*

Looking downwards

NCEA Level 8 Stats 2013

Let's drag ourselves out of the mire - put aside teething problems and lift up our eyes!







In December, IT service company EMC published its Data Scientist Study 2011, which forecast a shortage of people with suitable skills for at least the next five years.

The study – which EMC said was the largest-ever global survey of the data science community spanning the US, UK, France, Germany, India and China – comes at a time when the volume of data being generated is

growing exponentiall



1212121212121212

* For an organization to compete today, it has to embrace modern business complexity and harness the incredible power of data. A successful business now requires a proficient team of managers with both an in-depth knowledge of analytical tools and expertise in business.

> Dr. Yuri Medvedev Chief Mathematician Bank of Montreal

So what's driving all of this?



DEPARTMENT OF STATISTIC

So the data world ... is getting a whole lot bigger



So the data world ... is getting a whole lot bigger

- There is an explosion in the ...
 - quantities of data being collected
 - conceptions of what constitutes data
 - settings in which it can arise
 - ways of looking at it

Caught in the path of The Data Deluge

So here we are ...

What do

we do ??!!

"Stop"..." I wapt to get off !! "







The democratisation of data



DEPARTMENT OF STATISTICS

The democratisation of data

has 2 key ingredients:

- 1. Access to data
 - Open data movement …
- 2. Ability to understand it - for most people
 - » And that's where we come in!

"Power to the people!"

Competitive economies

NZAMT 201

With the data world ... getting a whole lot bigger

Can't just keep illuminating same small patch

- Need to get much ...
 - further
 - faster
 - & with better comprehension

Opening Up the World of data



VISION Statement for Early Statistics

To create excitement about
"What I can do with data &

What data can do for me"





Shift of focus

From:

• "What are all the things I have to do to get the output I'm meant to produce?"

To:

- "What are the questions?"
- "What can I see?"
- "What does that tell me?"



Using *"enabling software"* to "cut out the middle man"

Addressing the Need for Speed



 VIT-type software can facilitate a fast, accessible way in to understanding basic inferential conceptions

My "vision" is

initially create an appreciation of a very wide array of data types and what they can tell you

• and only then back fill the details (for those who need it)

Online Activities

From Census At School NZ

27. In the last seven days, which of these online activities have you done? (You may tick more than one)

Downloaded or listened online to music

- Downloaded or viewed online video (e.g. YouTube, TV shows, movies)
- Played online game(s)
- Exercise Kept in touch with friends (e.g. through instant messaging, Bebo or Facebook)
- Researched topics related to school work
- Have been online but did none of the above activities
- Have not been online in the last seven days



How do people deal with MR questions?

Pie charts that add to more than 100%!



Data Structure									
AX	AY	AZ	BA	BB	BC	BD	BE	BF	BG
technone	cellmonth	onlinemusic	onlinevideo	onlinegame	onlinefriend	onlineschool	onlineother	onlinenone	bedtime
no	12	yes	yes	no	yes	yes	no	no	11:45:00
no	1	yes	yes	no	yes	no	no	no	12:30:00
no	NULL	yes	yes	yes	yes	yes	no	no	10:00:00
no	36	yes	yes	no	yes	no	no	no	10:00:00
no	10	no	no	no	no	no	no	no	11:30:00
no	12	no	no	yes	yes	yes	no	no	10:00:00
no	1	no	yes	yes	yes	no	no	no	9:00:00
no	17	no	yes	yes	yes	yes	no	no	10:00:00
no	13	no	yes	yes	yes	yes	no	no	8:45:00
no	2	no	no	no	no)	8:45:00
no	3	no	yes	no	no	"T w	as bli	ind	10:00:00
no	13	yes	yes	no	yes				10:00:00
no	11	yes	yes	no	yes	707	no	no	10:30:00
yes	7	no	no	no	no	n	no	yes	8:00:00
yes	NULL	no	no	no	no	yes	no	no	12:00:00
no	12	yes	yes	yes	yes	yes	no	no	NULL
no	15	yes	yes	yes	yes	no	no	no	10:30:00
no	13	yes	yes	no	yes	yes	no	no	10:15:00
no	14	yes	no	yes	yes	yes	no	no	10:30:00
no	5	yes	yes	yes	yes	yes	no	no	10:30:00
no	5	no	no	yes	no	yes	no	no	8:00:00
no	8	yes	yes	no	yes	no	no	no	22:30:00
no	5	no	yes	yes	yes	no	no	no	10:30:00
no	NULL	no	no	yes	no	no	no	no	10:00:00
no	9	yes	yes	yes	no	yes	no	no	9:30:00
20	1/	VAC	VAC	VAC	VAC	VAC	20	200	11.00.00





Fundamental early-stat "big-data" question

- Of the ways in which we teach students to look at data ...
 - what still works?

and

• what needs to change? (as datasets get bigger)

Big statistical lessons

" Data beats anecdote"

– David Moore

NZAMT 201







Big statistical lessons

Data beats anecdote"

- David Moore

NZAMT 201

Big statistical lessons

Looking at the world using data is like looking through a window with ripples in the glass

"What I see ... is not quite the way it really is"







Big statistical lessons

The *pleasures of discovery* in data need to be accompanied by...

statistical "safe sex" messages

NIVERSITY OF AUCKLAND











Fact versus Artefact

Biggest difference in "detailed features"

- With smaller data sets, we only estimate & interpret gross/global features
 - Realise detailed features indistinguishable
 from randomness
 - Lucky by-product: deficiencies in data at a very detailed level may not hurt that much
 - "Artefact filtering" in terms of "significance", etc.
 - In statistical modelling, details swept up into "random error"
 - Lucky by-product: simple, comparatively understandable models
 » e.g. no high-order interactions

DEPARTMENT OF STATISTICS

Fact versus Artefact

Biggest difference in "detailed features"

With very big data sets ...

- many detailed features are not "random"
 - Large numbers of effects are "significant"
 - but may not be real either
- "Artefact detection" has to consider distorting filters that may be "sieving" the streams of data we get to see
 - cf. staring through a distorting camera lens

Fact versus Artefact

Biggest difference in "detailed features"

With very big data sets ...

- Quite detailed features are not random
 - Large numbers of effects are "significant"
 - That which was swept into "random error" becomes structure to be investigated/understood/interpreted
- "Artefact detection" has to consider distorting filters that may be "sieving" the streams of data we get to see
- Models likely have to smooth over smaller features we know are present in order to be understandable

DEPARTMENT OF STA

NZAMT 201

NZAMT 2013

Good news stories

- 1. Most of our graphs "scale"
- The huge conceptual leap comes ...
 - when move from representing individual points
 (concrete connections to data)
 - to depicting abstractions
 - Summaries (e.g. counts and proportions), densities, ...
- Once that hurdle has been jumped ...

size doesn't matter

- Well not much anyway!
 - "2nd order effects" that are relatively easy to understand

Good news stories

Even "individual-points" plots can be surprisingly useful with quite large data sets

- Will now use Mondrian ... (by Martin Theus)
 - to look at some quite large data sets using the humble scatterplot







10,000 data points NHANES 2003-4









What have we used?

- The humble scatter plot with
 - Varying levels of transparency
 - Varying point size
 - Zooming
- Very simple ideas that help us see a lot
- Being able to vary them dynamically is important

Earthquakes in Japan Since 1900

Ef Like 6 5 retwe

Author: Peter Aldhous published on New Scientist

Peter Aldhous takes us back in time with this viz of Japanese Earthquakes since 1900. Obviously, the island nati seismic activity.





NZAMT 201

NZAMT 201

THE UNIVERSITY OF AUCKLANE DEPARTMENT OF STATISTICS



Conclusions

- The world is changing & so must we
- The data world is exploding
 - As is its impact on our lives
 - Consequentially, the wider populace has to become better prepared to deal with it
 - "the democratisation of data"
 - Ability to understand as well as to access

Conclusions

• "the democratisation of data"

•

- Ability to understand as well as to access
- & who has the reach to make the biggest impact on the (widespread) ability to understand?
- So, who are the really key players here?

Will change in statistics ever stop?

- We stand where information technology meets human understanding
- So the answer must be

"Not if we are doing our jobs it won't"

- That said ...
 - we first need time ...
 - to get to grips with the changes already made
 - & to get the inevitable bugs out of our new systems
- Good news is that the new curriculum now allows ...
 - for on-going evolutionary advances
 - rather than widely-spaced step lurches

DEPARTMENT OF STATISTIC

NZAMT 2013





And find new & better ways to ...



to open up the data world for our students



Thank you