Some issues that are often overlooked in big data analytics

Limsoon Wong





What is big data and why

- Big data *a la* Gartner
 - Volume, velocity, variety
- Other
 characteristics
 - Veracity, v...

A practical definition "More than you know how to handle" • Why big data?

- Can collect cheaply, due to automation
- Can store cheaply, due to falling media prices
- Many success stories, where useful predictions were made with the data



- Much emphasis is on scaling issues
- But there are non-scaling-related issues that affect fundamental assumptions in current bioinformatics and statistical analysis
 - Big data may break analysis procedures in fundamental ways

Talk outline



4

- Forgotten assumptions
 - The 1st "I" in I.I.D.
 - The 2nd "I" in I.I.D.
- Overlooked information
 - Non-associations
 - Context
- More may not be better
 - Protein complexes
 - Causal genes



"All those in favor say 'Aye." "Aye." "Aye." "Aye." "Aye." "Aye." "Aye."

Forgotten assumptions

THE 1ST "I" IN I.I.D.

CIKM2014, Shanghai, 4 November 2014

	1	+ · c	entiles	R.
Degrees of the length of Array o [*] -100	Estimates in Ibs.	Observed deviates from 1907 lbs.	n Normal p.e =17	 Excess of Observed over Normal
5	1074	- 133	- 90	+43
10	1109	- 98	- 70	+ 28
15	1126	- 81	~ 57	+24
20	1148	· - 59	- 46	+13
1 25	1162	- 45	- 37	+ \$
30	1174	- 33	- 29	+ 4
35	1151	- 26	- 21	+ 5
40	1188	- 19	- 14	+ 5
45	1197	- to	- 7	+ 3
<i>m</i> 50	1207	0	0	· 0
55	1214	. + 7	+ 7	0
60	1219	+ 12	+14	- 2
65	1225	+ 18	. + 21	- 3
70	1230	+ 23	+ 29	- 0
23 75	1236	+ 29	+ 37	- 8
80	1243	+ 36	+ + 40i	- 10
85	1254	+ 47	+ 57	- 10
90	1267	+ 52	+70	- 18
95	1293	, + 86	+ 90	- 4

Distribution of the estimates of the dressed weight of a particular living ox, made by 787 different persons.

Vox Populi

Francis Galton, *Nature*, 75(1949):450-451, March 1907

"[The] middlemost estimate is 1207lb., and the weight of the dressed ox proved to be 1198 lb.; so the *vox populi* was in this case 9 lb., or 0.8 per cent. of the whole weight too high... This result is ... more creditable to the trustworthiness of a democratic judgment than might have been expected."

 q_1 , q_3 , the first and third quartiles, stand at 25° and 75° respectively. *m*, the median or middlemost value, stands at 50°. The dressed weight proved to be 1108 lbs.

Experiments on social influence

Lorenz et al., PNAS, 108(22):9020-9025, 2011



- 12 groups, 12 subjects each
- Each subject solves 6 different estimation tasks regarding geographical facts and crime statistics
- Each subject responds to 1st question on his own
- After all 12 group members made estimates, everyone gives another estimate, 5 consecutive times

- Different groups based their 2nd, 3rd, 4th, 5th estimates on
 - Aggregated info of others' from the previous round
 - Full info of others' estimates from all earlier rounds
 - Control, i.e. no info
- Two questions posed for each of the three treatments
- Each declares his confidence after the 1st and final estimates



Wisdom of the crowd

Table 1. The wisdom of crowd effect exists with respect to the geometric mean but not with respect to the arithmetic mean

		WISCONT-OT-CLOWC		
Question	True value	Arithmetic mean	Geometric mean	Median
1. Population density of Switzerland	184	2,644 (+1,337.2%)	132 (–28.1%)	130 (–29.3%)
2. Border length, Switzerland/Italy	734	1,959 (+166.9%)	338 (-54%)	300 (–59.1%)
3. New immigrants to Zurich	10,067	26,773 (+165.9%)	8,178 (–18.8%)	10,000 (-0.7%)
4. Murders, 2006, Switzerland	198	838 (+323.2%)	174 (–11.9%)	170 (–14.1%)
5. Rapes, 2006, Switzerland	639	1,017 (+59.1%)	285 (-55.4%)	250 (-60.9%)
6. Assaults, 2006, Switzerland	9,272	135,051 (+1,356.5%)	6,039 (–34.9%)	4,000 (–56.9%)

Wisdom-of-crowd aggregation

The aggregate measures arithmetic mean, geometric mean, and median are computed on the set of all first estimates regardless of the information condition. Values in parentheses are deviations from the true value as percentages.

- 1st estimates not normally distributed
- They are lognormally distributed

⇒ Subjects had problems choosing the right order of magnitude

Social influence effect



Social influence diminishes diversity in groups
 ⇒ Groups potentially get into "group think"!

of Singapore

Range reduction effect

aggregated

information



full

information

10

no information



- Group zooms into wrong estimate
- Truth may even be outside all estimates

Confidence effect



	(1) groups' wisdom-of- crowd indicator	(2) individuals' increase in confidence
intercept	3.92*** (15.1)	-0.049 (-0.97)
aggregated information	-1.39** (-3.29)	0.20** (3.08)
full information	-0.98* (-2.21)	0.28*** (4.14)
N	288	864

t statistics in parentheses (robust std. errors), * p < 0.05, ** p < 0.01, *** p < 0.001

Let $x_1 \dots X_n$ be the sorted estimates. Wisdom of the crowd indicator is max{i | $x_i \leq truth \leq x_{n-i+1}$ }

 Opinion convergence boosts individuals' confidence in their estimates despite lack of collective improvement in accuracy

CIKM2014, Shanghai, 4 November 2014

Copyright 2014 © Limsoon Wong

Social influence diminishes wisdom of the crowd



- Social influence triggers convergence of individual estimates
- The remaining diversity is so small that the correct value shifts from the center to the outer range of estimates
- ⇒ An expert group exposed to social influence may result in a set of predictions that does not even enclose the correct value any more!
- Conjecture: Negative effect of social influence is more severe for difficult questions

Related issue: People do not say what they really want to say





13

Stephen King, "Conflict between public and private opinion", *Long Range Planning,* 14(4):90-105, August 1981

"In fact, the evidence is very strong that there is a genuine difference between people's private opinions and their public opinions."



Forgotten assumptions

THE 2ND "I" IN I.I.D.

CIKM2014, Shanghai, 4 November 2014

Copyright 2014 © Limsoon Wong

Statistical tests



15

 Commonly used statistical tests (T-test, χ2 test, Wilcoxon rank-sum test, ...) all assume samples are drawn from independent identical distributions (I.I.D.)



How to ensure I.I.D.?

- In clinical testing, we carefully choose the sample to ensure I.I.D. so that the test is valid
 - Independent: Patients are not related
 - <u>Identical</u>: Similar # of male/female, young/old, ... in cases and controls

	А	В
lived	60	65
died	100	165

Note that sex, age, ... don't need to appear in the contingency table

- In big data analysis, and in many datamining works, people hardly ever do this!
 - Is this sound?



What is happening here?



Overall

	Α	В
lived	60	65
died	100	165

Looks like treatment A is better

Women

	Α	В
lived	40	15
died	20	5

History of heart disease

	Α	В
lived	10	5
died	70	50

Men

	Α	В
lived	20	50
died	80	160

No history of heart disease

	Α	В
lived	10	45
died	10	110

Looks like treatment B is better

Looks like treatment A is better

of Singapore

18

Sample not identically distributed



Overall

	Α	В
lived	60	65
died	100	165

Women

	Α	В
lived	40	15
died	20	5

History of heart disease

	Α	В
lived	10	5
died	70	50

Men

	А	В
lived	20	50
died	80	160

No history of heart disease

	Α	В
lived	10	45
died	10	110

Taking A

- Men = 100 (63%)
- Women = 60 (37%)
- Taking B
 - Men = 210 (91%)
 - Women = 20 (9%)

Men taking A

- History = 80 (80%)
- No history = 20 (20%)
- Men taking B
 - History = 55 (26%)
 - No history = 155 (74%)

Copyright 2014 © Limsoon Wong

Simpson's paradox in an Australian population census



19

Context	Comparing Groups	sup	P _{class=>50K}	p-value	
Race =White	Occupation = Craft-repair	3694	22.84%	1 00 10-19	
	Occupation = Adm-clerical		14.23%	1.00×10^{-13}	

Context	Extra attribute	Comparing Groups	sup	P _{class=>50K}
Race =White	Sex = Male	Occupation = Craft-repair	3524	23.5%
		Occupation = Adm-clerical	1038	24.2%
		Occupation = Craft-repair	107	8.8%
		Sex = Female	Occupation = Adm-clerical	2046

- Violation of the 2nd "I" of I.I.D.
- Btw, "men earn more than women" also violates the 2nd "I" in I.I.D.

Stratification



- Cannot test "H: Men earn more than women" directly because I.I.D. is violated
 - Different distributions of men & women wrt occupation
- Test instead
 - "S₁: For craftsmen, men earn more than women"
 - "S₂: For admin clerks, men earn more than women"
 -

where craftsmen, admin clerks, ... form an exhaustive list of disjoint occupations, provided each of $S_1, S_2, ...$ is valid

• Cf. Mantel test, Cochran test



Related issue: Sampling bias

"Dewey Defeats

Truman" was a famously incorrect banner headline on the front page of the Chicago Tribune on November 3, 1948, the day after incumbent United States President Harry S. Truman won an upset victory over Republican challenger and Governor of New York Thomas E. Dewey in the 1948 presidential election.



President-elect Truman holding the infamous issue of the *Chicago Tribune*, telling the press, "That ain't the way I heard it!"

The reason the Tribune was mistaken is that their editor trusted the results of a phone survey... Telephones were not yet widespread, and those who had them tended to be prosperous and have stable addresses.



Overlooked information

NON-ASSOCIATIONS

CIKM2014, Shanghai, 4 November 2014

Copyright 2014 © Limsoon Wong

We tend to ignore non-associations



23

- We have many technologies to look for associations and correlations
 - Frequent patterns
 - Association rules

- We tend to ignore non-associations
 - We think they are not interesting / informative
 - There are too many of them
- We also tend to ignore relationship between associations





• Dietary fat intake correlates with breast cancer

24



And like this...



Animal fat intake correlates with breast cancer

But not non-correlations like this.



Plant fat intake doesn't correlate with breast cancer

26

of Singapore



Yet there is much to be gained when we take both into our analysis

A: Dietary fat intake correlates with breast cancer

B: Animal fat intake correlates with breast cancer

C: Plant fat intake doesn't correlate with breast cancer ⇒ Given C, we can eliminate A from consideration, and focus on B!





Back to the Simpson's paradox

Context	Comparing Groups	sup	P _{class=>50K}	p-value	
Race =White	Occupation = Craft-repair	3694	22.84%	1.00 × 10 ⁻¹⁹	
	Occupation = Adm-clerical	3084	14.23%		

Context	Extra attribute	Comparing Groups	sup	P _{class=>50K}
Race =White	Sex = Male	Occupation = Craft-repair	3524	23.5%
		Occupation = Adm-clerical	1038	24.2%
		Occupation = Craft-repair	107	8.8%
		Sex = Female	Occupation = Adm-clerical	2046

- 2nd "I" in I.I.D. is violated
- Btw, "men earn more than women" also violates the 2nd "I" in I.I.D.

CIKM2014, Shanghai, 4 November 2014

It pays to look at relationship betw associations & non-associations

National Univer of Singapore 29

- A. Wrt craftsmen / admin clerks, there are more / less men than women
- B. Wrt men / women, craftsmen earn similar to admin clerks
- C. Wrt craftsmen / admin clerks, men earn more than women

 $P(m|c) > P(w|c) \Rightarrow P(m|c) > 50\%$

- $P(w|a) > P(m|a) \Rightarrow P(m|a) < 50\%$
- i.e. P(m| c) > P(m| a)
- P(\$ | m, c) ≈ P(\$ | m, a)
- $P($ | w, c) \approx P($ | w, a)$
- P(\$ | m, c) > P(\$ | w, c)
 - P(\$ | m, a) > P(\$ | w, a)

P(\$| c)

= P(\$, m| c) + P(\$, w| c)= P(\\$| m,c) P(m|c) + P(\\$|w,c) P(w|c) = [P(\\$|m,c) - P(\\$|w,c)] P(m|c) + P(\\$|w,c) > [P(\\$|m,a) - P(\\$|w,a)] P(m|a) + P(\\$|w,a) = P(\\$|m, a) P(m| a) + P(\\$|w, a) P(w| a) = P(\\$, m| a) + P(\\$, w| a) = P(\\$| a)

i.e., P(\$| c) ≥ P(\$| a)

"Craftsmen earn more than admin clerks" is an artefact

i.e., even if "craftsmen earn more than admin clerks" passes a valid statistical test, it is a derivative of A, B, C

context

/ˈkɒntɛkst/ Đ

noun

the circumstances that form the setting for an event, statement, or idea, and in terms of which it can be fully understood.

"the proposals need to be considered in the context of new European directives" synonyms: circumstances, conditions, surroundings, factors, state of affairs; More

 the parts of something written or spoken that immediately precede and follow a word or passage and clarify its meaning.

"skilled readers use context to construct meaning from words as they are read"

Overlooked information

CONTEXT

We tend to ignore context



31

- We have many technologies to look for associations and correlations
 - Frequent patterns
 - Association rules

- We tend to assume the same context for all patterns and set the same global threshold
 - This works for a focused dataset
 - But for big data where you union many things, this spells trouble

. . .



Formulation of a Hypothesis

- "For Chinese, is drug A better than drug B?"
- Three components of a hypothesis:
 - Context (under which the hypothesis is tested)
 - Race: Chinese
 - Comparing attribute
 - Drug: A or B
 - Target attribute/target value
 - Response: positive
- {{Race=Chinese}, Drug=A|B, Response=positive}



The right support threshold

{{Race=Chinese}, Drug=A|B, Response=positive}

Context	Comparing attribute	response= positive	response= negative
{Race=Chinese}	Drug=A	N ^A _{pos}	$N^A - N^A_{pos}$
	Drug=B	N ^B _{pos}	$N^B - N^B_{pos}$

- To test this hypothesis we need info:
 - N^A =support({Race=Chinese, Drug=A})
 - N^A_{pos} =support({Race=Chinese, Drug=A, Res=positive})
 - N^B =support({Race=Chinese, Drug=B})
 - N^B_{pos} =support({Race=Chinese, Drug=B, Res=positive})

⇒ Frequent pattern mining, but be careful with support threshold, need to relativize to context

CIKM2014, Shanghai, 4 November 2014

Relativizing to context



34

 Most people cannot set support threshold correctly when relativizing to context

A quick test!

- Suppose a test of a disease presents a rate of 5% false positives, and the disease strikes 1/1000 of the population
- Let's say people are tested randomly and a particular patient's test is positive
- What's the probability that he is stricken with the disease?

Answer



35

- P(d) = 0.1%
- P(pos| ~d) = 5%
- P(pos| d) = 100%, assuming 100% sensitivity
- P(pos) = P(pos| d) P(d) + P(pos| ~d) P(~d) ≈ 5%
- P(d| pos) = P(pos| d) P(d) / P(pos) = 0.1% / 5% = 2%
- I.e., the answer is 2%
- Did you guess 95% as the answer?



The right context

{{Race=Chinese}, Drug=A|B, Response=positive}

Context	Comparing attribute	response= positive	response= negative
{Race=Chinese}	Drug=A	N ^A _{pos}	$N^A - N^A_{pos}$
	Drug=B	N ^B _{pos}	$N^B - N^B_{pos}$

- If A/B treat the same single disease, this is ok
- If B treats two diseases, this is not sensible
- The disease has to go into the context



More may not be better

PROTEIN COMPLEXES

CIKM2014, Shanghai, 4 November 2014

Copyright 2014 © Limsoon Wong



Copyright 2014 © Limsoon Wong

Difficulties



39

- Cytochrome BC1 complex
 - Involved in electrontransport chain in mitochondrial inner membrane



Figure 1 PPI subgraph of the mitochondrial cytochrome bc1 complex. Nineteen interactions were detected between the ten proteins from the complex, while many extraneous interactions were detected. Five example proteins from transient interactions are down: NAB2 and UBI4 are involved in mRNA polyadenylation and protein ubiquitination, while PET9, SHY1, and COX1 are mitochondrial membrane proteins that are also involved in the electron-transport chain. The extraneous interactions around the complex makes its discovery difficult. All such network figures were generated by Cytoscape [30].

- Discovery of BC1 from
 PPI data is difficult
 - Sparseness of its PPI subnetwork
 - Only 19 out of 45 possible interactions were detected between the complex's proteins
 - Extraneous interactions with other proteins outside the complex
 - E.g., UBI4 is involved in protein ubiquitination, and binds to many proteins to perform its function

Copyright 2014 © Limsoon Wong

CIKM2014, Shanghai, 4 November 2014



Perhaps "big data" can help?

Composite network

 Vertices represent proteins, edges represent relationships between proteins. Put an edge betw proteins u, v, iff u and v are related according to any of the data sources

Data sourc	e	Databa	ase		Scoring me	thod
PPI		BioGR	ID, IntACT, MI	NT	Iterative Adju	stCD.
L2-PPI (indi	rect PPI)	BioGR	ID, IntACT, M	NT	Iterative AdjustCD	
Functional a	association	STRIN	G		STRING	3
Literature c	o-occurrence	PubMe	ed		Jaccard coef	ficient
		Yeast			Human	
	# Pairs	% co-complex	coverage	# Pairs	% co-complex	coverage
PPI	106328	5.8 %	55 %	48098	10%	14%
L2-PPI	181175	1.1%	18%	131705	5.5%	20%
STRING	175712	5.7%	89%	311435	3.1%	27%
PubMed	161213	4.9%	70%	91751	4.3%	11%
All	531800	2.1 %	98 %	522668	3.4%	49%

CIKM2014, Shanghai, 4 November 2014

Copyright 2014 © Limsoon Wong

Yong, et al. Supervised maximum-likelihood weighting of composite protein networks for complex prediction. *BMC Systems Biology*, 6(Suppl 2):S13, 2012



41

More is not always better, unless.



SWC-weighted network



While proteins in BC1 become fully connected in the composite network, there is a blow-up in extraneous proteins. So clustering won't discover the complex, unless you know how to remove the extraneous proteins



More may not be better

CAUSAL GENES

NUS National University of Singapore

Gene expression analysis challenge

- Low % of overlapping genes from diff expt in general
 - Prostate cancer
 - Lapointe et al, 2004
 - Singh et al, 2002
 - Lung cancer
 - Garber et al, 2001
 - Bhattacharjee et al, 2001
 - DMD
 - Haslett et al, 2002
 - Pescatori et al, 2007

Datasets	DEG	POG
Prostate	Тор 10	0.30
Cancer	Тор 50	0.14
	Top100	0.15
Lung	Тор 10	0.00
Cancer	Тор 50	0.20
	Top100	0.31
DMD	Тор 10	0.20
DMD	Тор 50	0.42
	Top100	0.54

Zhang et al, Bioinformatics, 2009

Copyright 2014 © Limsoon Wong



Biology to the rescue?



- Each disease phenotype has some underlying cause
- There is some unifying biological theme for genes that are truly associated with a disease subtype
- Uncertainty in selected genes can be reduced by considering biological processes of the genes
- The unifying biological theme is basis for inferring the underlying cause of disease subtype

Database	Remarks 🔤 🛞
KEGG	KEGG (http://www.genome.jp/kegg) is one of the best known pathway databases (Kanehisa <i>et al.</i> , 2010). It consists of 16 main databases, comprising different levels of biological infor- mation such as systems, genomic, etc. The data files are down- loadable in XML format. At time of writing it has 392 path- ways.
WikiPathways	WikiPathways (http://www.wikipathways.org) is a Wikipedia-based collaborative effort among various Big data o labs (Kelder et al., 2009). It has 1,627 pathways of which 369 are human. The content is downloadable in GPML format.
Reactome	Reactome (http:://www.reactome.org) is also a collaborative effort like WikiPathways (Vastrik <i>et al.</i> , 2007). It is one of the largest datasets, with over 4,166 human reactions organized into 1,131 pathways by December 2010. Reactome can be down- loaded in BioPax and SBML among other formats.
Pathway Commons	Pathway Commons (http://www.pathwaycommons.com) col- lects information from various databases but does not unify the data (Cerami <i>et al.</i> , 2006). It contains 1,573 pathway 564 organisms. The data is returned in BioPax format
PathwayAPI	PathwayAPI (http://www.pathwayapi.com) contains unified human pathways obtained from a merge of WikiPathways and Ingenuity® Knowledge Base (Sol 2010). Data is downloadable as a SQL dump or as a and is also interfaceable in JSON format.

Goh, et al. *Proteomics*, 12(4-5):550-563, 2012.

Copyright 2014 © Limsoon Wong

Soh et al. Consistency, Comprehensiveness, and Compatibility of Pathway Databases. BMC Bioinformatics, 11:449, 2010.



45



Overlap Analysis: ORA



ORA tests whether a pathway is significant by intersecting the genes in the pathway with a pre-determined list of DE genes (we use all genes whose t-statistic meets the 5% significance threshold), and checking the significance of the size of the intersection using the hypergeometric test

S Draghici et al. "Global functional profiling of gene expression". *Genomics*, 81(2):98-104, 2003.



Disappointing Performance

upregulated in DMD





Issue #1 with ORA

- Its null hypothesis basically says "Genes in the given pathway behaves no differently from randomly chosen gene sets of the same size"
- This null hypothesis is obviously false
- \Rightarrow Lots of false positives



 A biological pathway is a series of actions among molecules in a cell that leads to a certain product or a change in a cell. Thus necessarily the behavour of genes in a pathway is more coordinated than random ones



Issue #2 with ORA

- It relies on a predetermined list of DE genes
- This list is sensitive to the test statistic used and to the significance threshold used
- This list is unstable regardless of the threshold used when sample size is small



t-test p.value(s)



Issue #3 with ORA

- It tests whether the entire pathway is significantly differentially expressed
- If only a branch of the pathway is relevant to the phenotypes, the noise from the large irrelevant part of the pathways can dilute the signal from that branch





performance improves

upregulated in DMD



51



What have we learned?

- More data can offer a more complete picture, fill in gaps, etc.
- More data can also introduce noise into an analysis
- Unless you know how to tame this noise, more data may not lead to a better analysis

- Mechanical application of statistical and data mining techniques often does not work
- Must understand statistical and data mining tools & the problem domain
 - Must know how to logically exploit both