CS4220 Knowledge Discovery Methods for Bioinformatics Unit1: Essence of Knowledge Discovery (Part B: Art of Statistical Analysis)

Wong Limsoon



### Outline



2

### Forgotten assumptions

- Normal distribution
- I.I.D.
- Proper design of experiment
- Domain-specific laws
- Overlooked information
  - Non-associations
  - Context



3



#### **Forgotten assumptions**

## **NORMAL DISTRIBUTION**



Copyright 2018 © Limsoon Wong

### Wisdom of the crowd Lorenz et al., PNAS, 108(22):9020-9025, 2011



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

			auggregation		
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 addregation

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.

- Estimates not normally distributed
- They are lognormally distributed

⇒ Subjects had problems choosing the right order of magnitude

### Time for Exercise #1



5

- Suppose you are given a set S of values (e.g. the age of a group of people). Choose a number or value x so that x would be a good representative of the values in S when
  - S is normally distributed
  - S is log-normally distributed
  - S has some arbitrary distribution
- What is the general principle underlying your choices?



6

# Me: I'm finally happy. Life: Lol, wait a sec.

and what held yesterday may not hold today



### 2007 Financial Crisis





 All of them religiously check VaR (Value at Risk) everyday

- VaR measures the expected loss over a horizon assuming normality
- "When you realize that VaR is using tame historical data to model a wildly different environment, the total losses of Bear Stearns' hedge funds become easier to understand. It's like the historic data only has rainstorms and then a tornado hits." – New York Times, 2 Jan 2009
- You can still turn things into your advantage if you are alert: When VaR numbers start to miss, either there is something wrong with the way VaR is being calculated, or the market is no longer normal



# Forgotten assumptions

Copyright 2018 © Limsoon Wong



### 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 1<sup>st</sup> question on his own
- After all 12 group members made estimates, everyone gives another estimate, 5 consecutive times

- Different groups based their 2<sup>nd</sup>, 3<sup>rd</sup>, 4<sup>th</sup>, 5<sup>th</sup> 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 1<sup>st</sup> and final estimates

### Social influence effect



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

10

of Singapore

### Range reduction effect

aggregated

information



full

information

#### no information



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

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



CS4220



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**

## **PROPER DESIGN OF EXPT**

Copyright 2018 © Limsoon Wong

### Design of experiments



15

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

	A	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

	A	В
lived	60	65
died	100	165

#### Men

	A	В
lived	20	50
died	80	160

#### History of heart disease

Α

40

20

Women

lived

died

CS4220

	Α	В
lived	10	5
died	70	50

В

15

5

#### No history of heart disease

	A	В
lived	10	45
died	10	110

#### Looks like treatment A is better

#### Looks like treatment B is better

#### Looks like treatment A is better



16

# A/B sample not identical in other attributes





#### Overall

	Α	В
lived	60	65
died	100	165

#### Women

	Α	В
lived	40	15
died	20	5

#### History of heart disease

	А	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%)



# Simpson's paradox in an Australian population census



18

Context	Comparing Groups	sup	P <sub>class=&gt;50K</sub>	p-value	
Race =White	Occupation = Craft-repair	3694	22.84%	1 00 10-19	
	Occupation = Adm-clerical	3084	14.23%	1.00 × 10 <sup>18</sup>	

Context	Extra attribute	Comparing Groups	sup	P <sub>class=&gt;50K</sub>
Race =White	Sex = Male	Occupation = Craft-repair	3524	23.5%
		Occupation = Adm-clerical	1038	24.2%
	Sex = Female	Occupation = Craft-repair	107	8.8%
		Occupation = Adm-clerical	2046	9.2%

 Craft-repair/Adm-clerical sample not identical in other aspects



### Time for Exercise #2



19

18

 Slide #18 suggests that men earn more than women. How would you verify this hypothesis? Should you do a chi-square test using the table shown below?

### Simpson's paradox in an Australian population census

Context	Comparing Groups		sup	P <sub>class=&gt;50K</sub>			p-value
Raco -W/bito	Occupation = Craft-repair		3694	22.84%		1.0010-19	
Race - White	Occupation = Adm-clerical		3084	14.2	4.23%		.00 × 10 ···
Context	Extra attribute	Comparing Groups			sup		P <sub>class=&gt;50K</sub>
Race =White	Sex = Male	Occupation = Craft-repair			352	4	23.5%
		Occupation = Adm-clerical			103	8	24.2%
	Sex = Female	Occupation = Craft-repair			107	7	8.8%
		Occupation = Adm-clerical			204	6	9.2%

Craft-repair/Adm-clerical sample not identical in other aspects

Copyright 2016 © Limsoon Wong

	Earn <50k	Earn ≥50k
Sex = Male	3483 (76%)	1079 (24%)
Sex = Female	1955 (91%)	198 (9%)

### 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.



**Forgotten assumptions** 

## **DOMAIN-SPECIFIC LAWS**

Copyright 2018 © Limsoon Wong







			0	Group			
SNP	Genotypes	Contr	ols [n(%)]	Cases	s [n(%)]	χ²	P value
rs???	AA	1	0.9%	0	0.0%		4.78E-21 <sup>b</sup>
	AG	38	35.2%	79	97.5%		
	GG	69	63.9%	2	2.5%		

Abbreviation: SNP, single nucleotide polymorphism.

24

### Time for Exercise #3



25

 Slide #24 says the contingency table looks suspicious. Why?





Copyright 2018 © Limsoon Wong



#### **Overlooked information**

## **NON-ASSOCIATIONS**

Copyright 2018 © Limsoon Wong

### We tend to ignore non-associations



28

- 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

29



30

### 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

of Singapore



32

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!





#### 33

### 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



34

- 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



35

- "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}



36

### The right support threshold

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

Context	Comparing attribute	response= positive	response= negative	
{Race=Chinese}	Drug=A	N <sup>A</sup> <sub>pos</sub>	$N^A - N^A_{pos}$	
	Drug=B	N <sup>B</sup> <sub>pos</sub>	$N^B - N^B_{pos}$	

- To test this hypothesis we need info:
  - N<sup>A</sup> =support({Race=Chinese, Drug=A})
  - N<sup>A</sup><sub>pos</sub> =support({Race=Chinese, Drug=A, Res=positive})
  - N<sup>B</sup> =support({Race=Chinese, Drug=B})
  - N<sup>B</sup><sub>pos</sub> =support({Race=Chinese, Drug=B, Res=positive})

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

### The right context



37

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

Context	Comparing attribute	response= positive	response= negative
{Race=Chinese}	Drug=A	N <sup>A</sup> <sub>pos</sub>	$N^A - N^A_{pos}$
	Drug=B	N <sup>B</sup> <sub>pos</sub>	$N^B - N^B_{pos}$

- 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

### Time for Exercise #4



38

 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 is the probability that he is stricken with the disease?



40

# What have we learned?

- 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





Undamaged plane (left). A plane shaded everywhere bullets struck returning aircraft (right).

- "It is so easy to make bad inferences with data... there's a creative part of understanding quantitative data that requires a sort of artistic or creative approach to research."
  ---Nate Bolt
- http://www.fastcodesign.com/1671172/how-a-story-from-world-war-ii-shapes-facebook-today

41