

Connecting technology to intelligence

www.scalnyx.com

Simpson's Paradox analyzed through Causal Reasoning

pyAgrum User day 15/06/2021

Santiago Cortijo, PhD, Head R&D santiago.cortijo@scalnyx.com

SCALNYX S.A.S. 5 Avenue Ingres, 75016 Paris, France info@scalnyx.com - Tel : +33 6 45 32 83 70



This presentation is based on:

- Understanding the Simpson's Paradox, by Judea Pearl (2013)
- 2. Causality Course, by Brady Neal:

https://www.bradyneal.com/causal-inference-course



Correlation *≠* **Causation**

Motivating examples



Correlation \neq **Causation (1)**



Sleep with shoes



Headache next morning.



Correlation \neq **Causation (2)**



Sleep with shoes

Headache next morning.



Correlation \neq **Causation (3)**





SCALNYX BY DEBIGN





Randomized Control Trials (RCT), (2)



RCT eliminates confounding association

Х



Randomized Control Trials (RCT), (3)



RCT eliminates confounding association

Sadly, RCT can are often expensive, unethical or even impossible.



Observational studies

How do we measure causal effects



Motivating example : Treatment to choose (1)

T: Treatment (A or B)



Death rate table:

Т	Total
A	16 % (240 / 1500)
В	19 % (105 / 550)



Decision Makers : Need to invest in a Treatment. Which one? which is better, A or B ?



Motivating example : Treatment to choose (2)

T: Treatment (A or B)

Death rate table, with mode details about symptoms intensity:

Т	Mild Symptoms.	Severe Symptoms	Total
А	15 %	30 %	16 %
	(210 / 1400)	(30 / 100)	(240 / 1500)
В	10 %	20 %	19 %
	(5 / 50)	(100 / 500)	(105 / 550)



Decision Makers: Wait... what !?



Motivating example : Treatment to choose (3)

T: Treatment (A or B)

Т	Mild	Severe	
A	15 % (210 / <u>1400</u>)	30 % (30 / <u>100</u>)	16 % (240 / 1500)
В	10 % (5 / 50)	20 % (100 / 500)	19 % (105 / 550)

$\frac{1400}{1500}$ 0.15 +	$\frac{100}{1500}$ 0.30 = 0.16
$\frac{50}{550}$ 0.10 +	$rac{500}{550} 0.20 = 0.19$



Motivating example : When is B prefered ?

Causal Graph:

Symptom intensity (C) causes Treatment (T)





Mild Symptoms



Severe Symptoms





Treatment A



Treatment B



Motivating example : When is A prefered ?

Causal Graph:

Treatment (T) causes Symptom intensity (C), due to waiting list



Treatment ATreatment BImage: Constrained a state of the state of t



Measure causal effects in observational data (1)

For our example, we adjust/control for confounders



 $\mathrm{E}[Y|t,C=0] \hspace{0.2cm} \mathrm{E}[Y|t,C=1] \hspace{0.2cm} \mathrm{E}[Y|t]$



$\mathrm{E}[Y \mid$	do(T	= t)] =	$\mathrm{E}_C\mathrm{E}[Y]$	$\mid T =$	t = t]
---------------------	------	---------	-----------------------------	------------	-----------

SCALNYX BY DEBIGN

Т	Mild	Severe	0,0	causal	
A	15 % (210 / 1400)	30 % (30 / 100)	16 % (240 / 1500)	19.4 %	$rac{1450}{2050} 0.15 + rac{600}{2050} 0.30 = 0.194$
В	10 % (5 / 50)	20 % (100 / 500)	19 % (105 / 550)	12.9 %	$rac{1450}{2050}0.10 + rac{600}{2050}0.20 = 0.129$

 $\mathrm{E}[Y|t,C=0] \hspace{0.1 in} \mathrm{E}[Y|t,C=1] \hspace{0.1 in} \mathrm{E}[Y|t] \hspace{0.1 in} \mathrm{E}[Y|do(t)]$

SCALNYX BY DEBIGN

Done with pyAgrum (1)

Generating data

```
[2]: filename1 = 'data1.csv'
columns = ['Treatment', 'Symptoms', 'Outcome']
data = \
        [['A', 'mild', 'dead']]*210 + \
        [['A', 'mild', 'alive']]*(1400-210) + \
        [['A', 'severe', 'dead']]*30 + \
        [['A', 'severe', 'dead']]*30 + \
        [['B', 'mild', 'dead']]*5 + \
        [['B', 'mild', 'alive']]*(50-5) + \
        [['B', 'severe', 'dead']]*100 + \
        [['B', 'severe', 'alive']]*(500-100)
shuffle(data)
df = pd.DataFrame(data=data, columns=columns)
```

df.to csv(filename1, index=False)

df

0	А	mild	alive
1	А	mild	alive
2	А	mild	alive
3	А	mild	alive
4	А	mild	alive
2045	А	mild	alive
2046	В	severe	alive
2047	А	mild	alive
2048	В	mild	alive
2049	А	mild	dead

Treatment Symptoms Outcome

2050 rows × 3 columns

[2]:

Done with pyAgrum (2) : Learning

Learning model from data (using constraints)



SCALNYX BY DEBION

Done with pyAgrum (3): Inference

gnb.flow.row(bn.cpt('Treatment'), bn.cpt('Symptoms'), bn.cpt('Outcome'))

							Outo	ome
					Treatment	Symptoms	alive	dead
	Treat	ment				mild	0.8499	0.1501
Symptoms	Α	в	Sym	Symptoms		severe	0.6995	0.3005
mild	0.9654	0.0346	mild	severe		mild	0.8980	0.1020
severe	0.1669	0.8331	0.7072	0.2928	в	severe	0.7999	0.2001

ie=gum.LazyPropagation(bn)
ie.evidenceImpact("Outcome",[])

Outcome				
alive dead				
0.8315	0.1685			

ie=gum.LazyPropagation(bn)
ie.evidenceImpact("Outcome",["Treatment"])

	Outcome			
Treatment	ent alive de			
А	0.8399	0.1601		
В	0.8088	0.1912		

Look, that Simspon's paradox still there !

Done with pyAgrum (4)

Generate (Draw) data from the model itself

<pre>filename2 = 'data2.csv'</pre>		Treatment	Symptoms	Outcome
<pre>dbgen = gum.BNDatabaseGenerator(bn)</pre>		ireaunent	Symptoms	outcome
dbgen.drawSamples(2050)	0	А	mila	alive
abgen.tocsv(Tilename2)	1	В	severe	alive
df2 = pd.read csv(filename2)	2	В	severe	alive
df2	3	В	severe	alive
	4	В	severe	alive
	2045	В	severe	dead
	2046	А	mild	alive
	2047	В	severe	alive
	2048	А	mild	alive
	2049	Α	mild	alive

2050 rows × 3 columns

Done with pyAgrum (5)

df2 = pd.read_csv(filename2)
ct = pd.crosstab(df2.Treatment, [df2.Outcome])
p_dead = ct['dead'] / (ct['alive'] + ct['dead'])
pd.DataFrame(p_dead.rename('prob(dead)'))

prob(dead)

Treatment

- A 0.153428
- **B** 0.164645

Simspon's paradox still there !

```
[15]: df2 = pd.read_csv(filename2)
  ct = pd.crosstab(df2.Treatment, [df2.Symptoms, df2.Outcome])
  p_dead_mild = ct[('mild', 'dead')] / ( ct[('mild', 'alive')] + ct[('mild', 'dead')])
  p_dead_severe = ct[('severe', 'dead')] /( ct[('severe', 'alive')] + ct[('severe', 'dead')] )
  pd.concat([p dead mild.rename('mild'), p dead severe.rename('severe')], axis=1).rename axis('Prob(dead)', axis=1)
```

[15]: Prob(dead) mild severe

Treatment

- A 0.140671 0.326733
- B 0.074074 0.173996



Done with pyAgrum (6)





Mild Symptoms

Severe Symptoms

learner=gum.BNLearner(filename1)
learner.addMandatoryArc('Symptoms', 'Treatment')
learner.addMandatoryArc('Symptoms', 'Outcome')
learner.addMandatoryArc('Treatment', 'Outcome')

bn = learner.learnBN()
bn

```
d1 = csl.CausalModel(bn)
cslnb.showCausalImpact(d1, "Outcome", doing="Treatment",values={})
```







Done with pyAgrum (8)



SOME CAUSAL DIAGRAMS (LIKE THESE) WON'T EVER GENERATE SIMPSON'S PARADOXES



Done with pyAgrum (9) : Backdoors



IF THERE IS A BACKDOOR, YOU MUST CONDITION ON IT.

IF THERE IS A BACKDOOR CHOOSE B

	Outcome			
Treatment	alive	dead		
A	0.8059	0.1941		
В	0.8693	0.1307		

Causal Model

Explanation : backdoor ['Symptoms'] found.



Done with pyAgrum (10): Backdoors

learner=gum.BNLearner(filename1)
learner.addMandatoryArc('Symptoms', 'Treatment')
learner.addForbiddenArc('Symptoms', 'Outcome')
learner.addMandatoryArc('Treatment', 'Outcome')|
bn6 = learner.learnBN()
d6 = csl.CausalModel(bn6, [("L1", ["Treatment", "Outcome"])], keepArcs=True)
cslnb.showCausalImpact(d6, "Outcome", doing="Treatment", values={})





Done with pyAgrum (11): Backdoors





IF THERE IS NO BACKDOOR, CHOOSE THE AGGREGATE.

 $P(Outcome | \hookrightarrow Treatment) = P(Outcome | Treatment)$

Causal Model

Explanation : Do-calculus computations



Summary

When facing the Simpson's Paradox:

- 1. Draw a causal Diagram
- 2. If the is a backdoor
 - a. If it is observed, the right decision it to use it as conditioning variableb. If it is unobserved, no conclusion can be made.
- 3. If there is no backdoor, the right decision is to use the aggregate result.





Thank you.