

Week 10

Live Discussion Session

Starts at 2.05pm

Message from Director of Teaching:
module evaluation questionnaire for Sem B taught modules is now open

- Important to complete the questionnaires and provide constructive feedback (including positive) in the free text areas.
- “Do not sit on the fence” when it comes down to selecting numerical scores in the evaluation. A neutral “3” in a 5-point scale is treated as a negative response. Go for a clearly positive 4-5 or a negative 1-2 as it sends a clearer message.

Question 3

- a) Pearl used a “ladder of causation” to explain why purely data-driven algorithms and classical statistics cannot achieve “true artificial intelligence”. State the simple words used by Pearl to describe the three rungs of the ladder and where on the ladder purely data-driven algorithms and classical statistics can reach. **[4 marks]**
- b) By using a simple medical example, explain the terms “association”, “intervention” and “counterfactual” used to describe different types of reasoning and how they relate to Pearl’s ladder of causation. **[6 marks]**
- c) A University has collected data on a large number of its students to determine whether the amount that a student spends on books influences their final degree performance. The data suggests that increased spending does improve performance. However, the data also shows that students who went to grammar or private schools tend to spend more on books and also achieve better degrees.

The observational data enables them to build a Bayesian network model as shown in Figure 1. There is a new proposal to buy £1000 worth of books for each student. How would you use the model to determine whether such an intervention would reduce the number of students failing their degree. **[6 marks]**

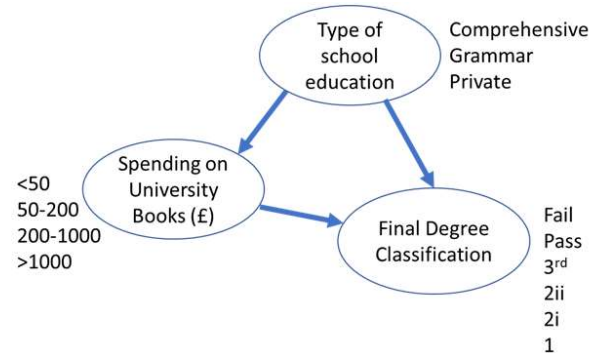


Figure 1 Student performance Bayesian network (with node states shown alongside nodes)

- d) A particular student Glenda achieved a 2i at this university. We know that Glenda spent less only £30 on books but we do not know what type of school she was educated in. How would we use the model to answer the question “Would Glenda have achieved a 1st class degree if she had spent at least £1000 on books instead of just £30”? **[9 marks]**

a) Rung 1: “Seeing” Rung 2: “Doing” Rung 3 “Imagining”

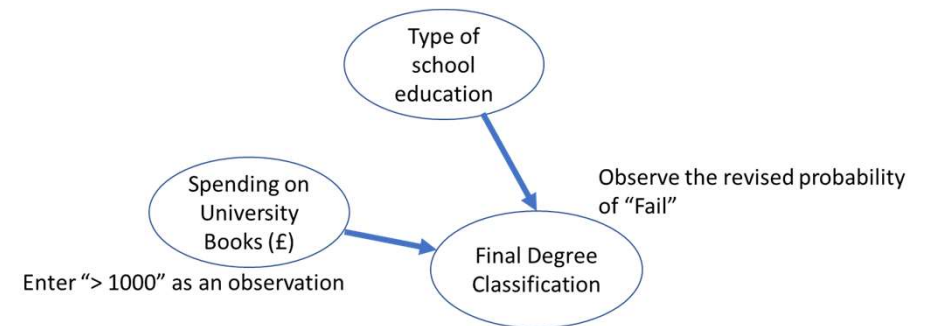
Purely data-driven algorithms and classical statistics can reach only “Seeing” rung 1

b) Association: Does use of this drug lead to improved recovery rates

Intervention: If I use this drug what is the probability I will recover.

Counterfactual: If I use the drug and recover – would I *still* have recovered if I had not used the drug.

c) First we need to note the marginal probability of ‘fail’ in the original model. Then we need to cut the link from ‘Type of school’ to ‘spending’:

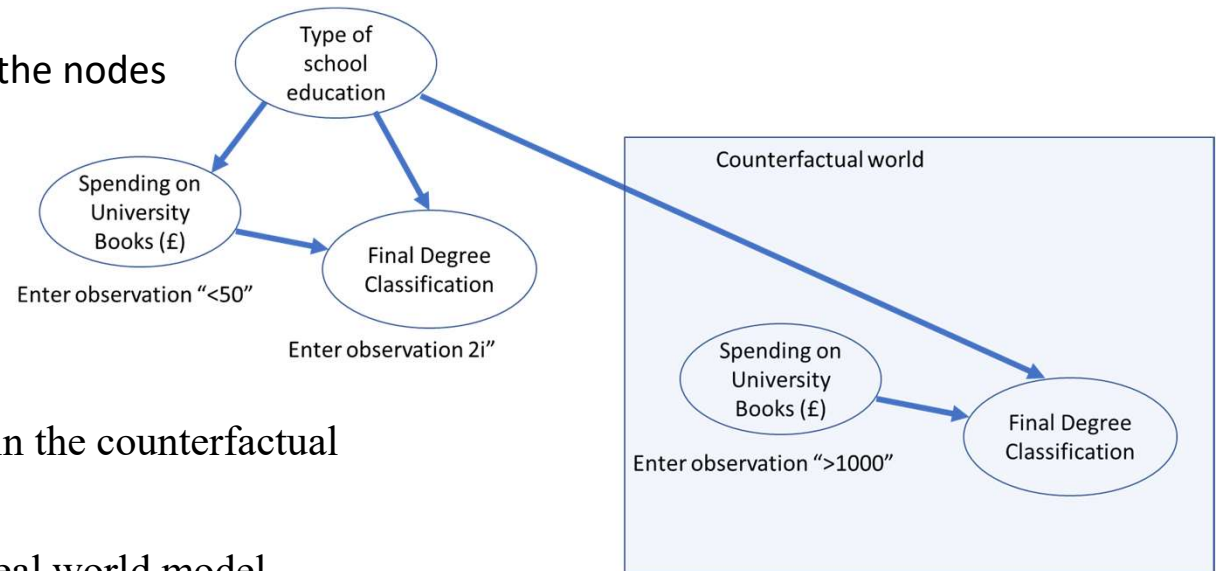


In this revised model enter “>1000” as an observation and run the model.

Observe the revised probability of fail. If it is less than the previous marginal then the intervention would be a success.

- d) A particular student Glenda achieved a 2i at this university. We know that Glenda spent less only £30 on books but we do not know what type of school she was educated in. How would we use the model to answer the question “Would Glenda have achieved a 1st class degree if she had spent at least £1000 on books instead of just £30”?
- [9 marks]**

Create the ‘twin network’ model by copying the nodes Spending and Final degree:



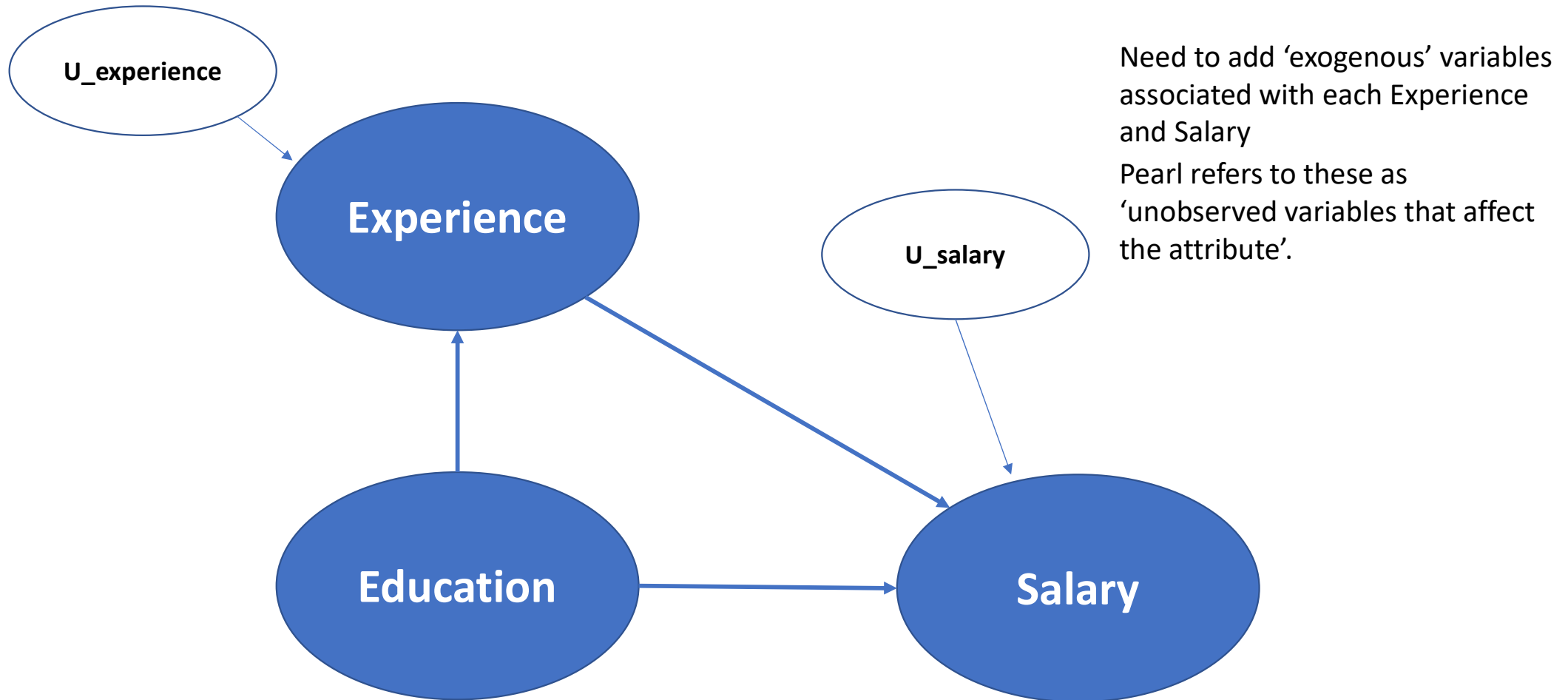
Cut the link from Type of School to Spending in the counterfactual world.

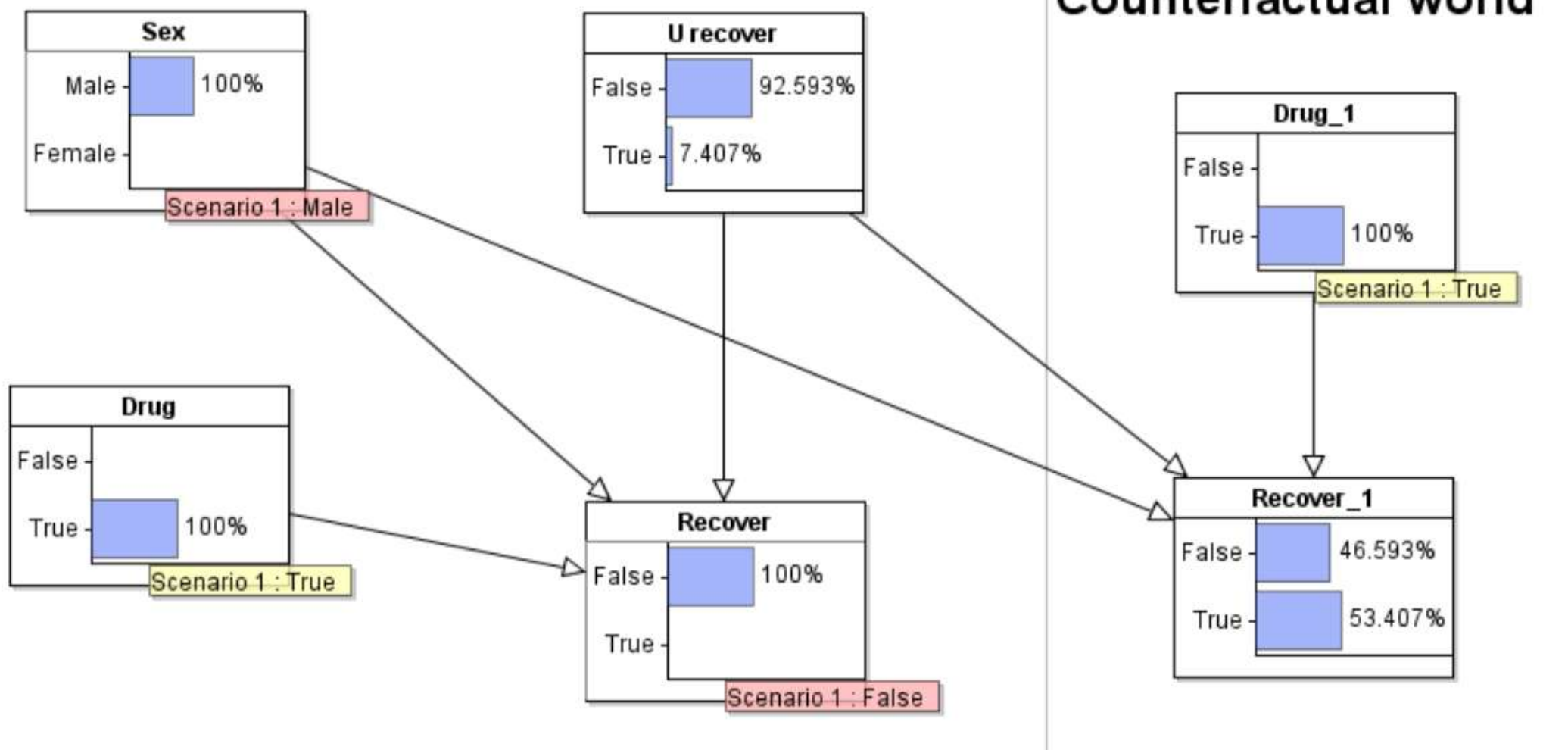
Enter the observations “<50” and “2i” in the real world model.

Run the model – this will update the “Type of School”

Enter the observation “>1000” in the counterfactual world and run the model again. Look at the probability of “1” in the final degree of the counterfactual world.

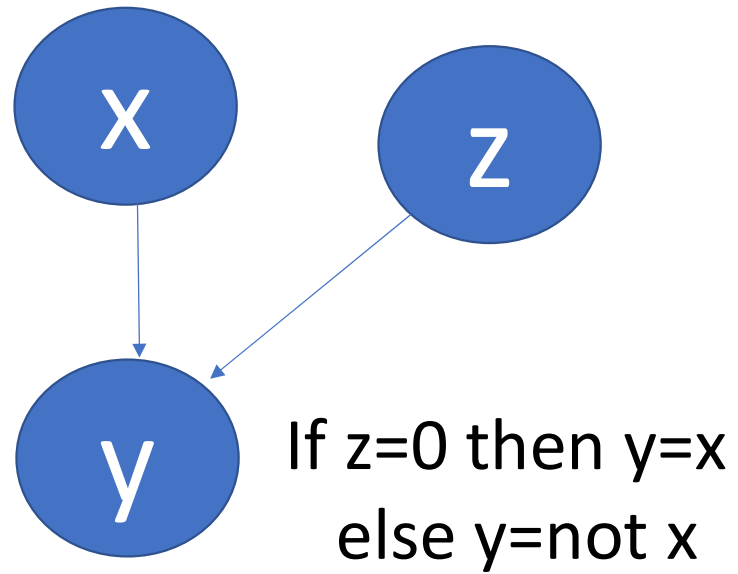
Pearl's example: data for potential outcomes





x	Outcome y
1	1
0	0
1	0
0	1
1	0
0	0
1	1
1	0
0	1
0	0
0	1
1	0
1	1
0	1
1	0
1	1
0	0
0	1

causation without correlation



Handling missing data values

The easily implemented methods assume data is missing at random

But whether data are collected from individuals by choice (e.g. online surveys, volunteers in studies, experiments) or automatically, there are almost inevitable systemic reasons accounting for missing data

Examples:

In any dataset where individuals are asked to give their salary, missing values are much more likely to be from those with very low or very high salaries

Males are systematically less likely to answer certain types of questions than females and vice versa

Any system automatically mining data about people from what is available online will have far more missing attribute values for older people