



# DATA ANALYTICS

## Unit 5: Advanced Techniques

---

**Swati Pratap Jagdale**

Department of Computer Science and  
Engineering

# DATA ANALYTICS

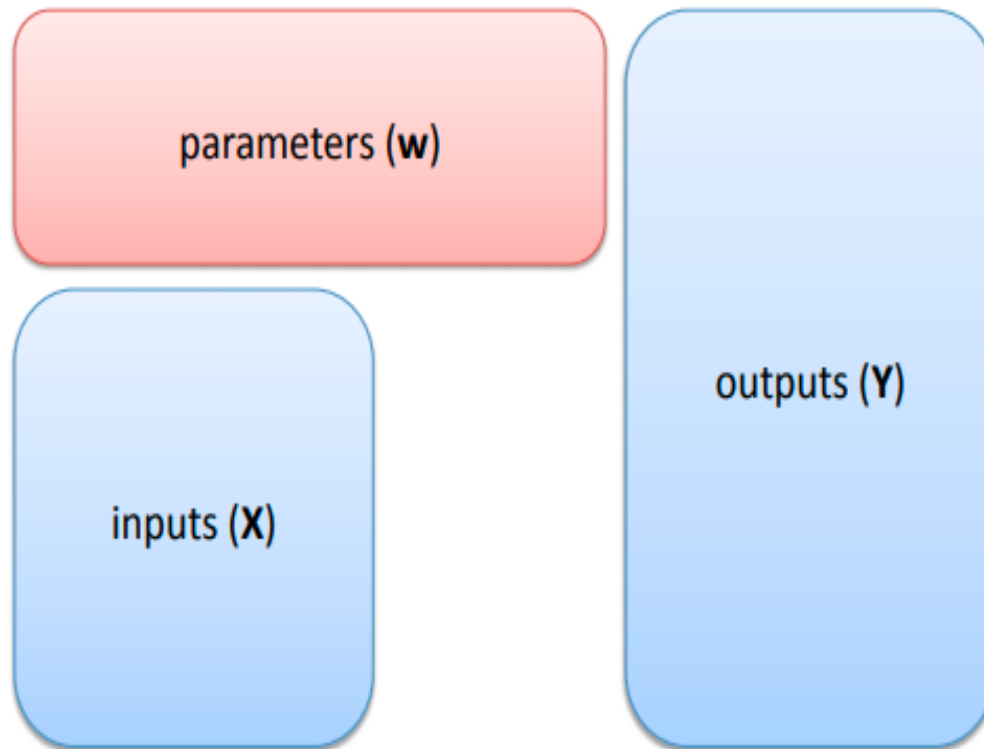
---

## Unit 5: Concept of hidden variables

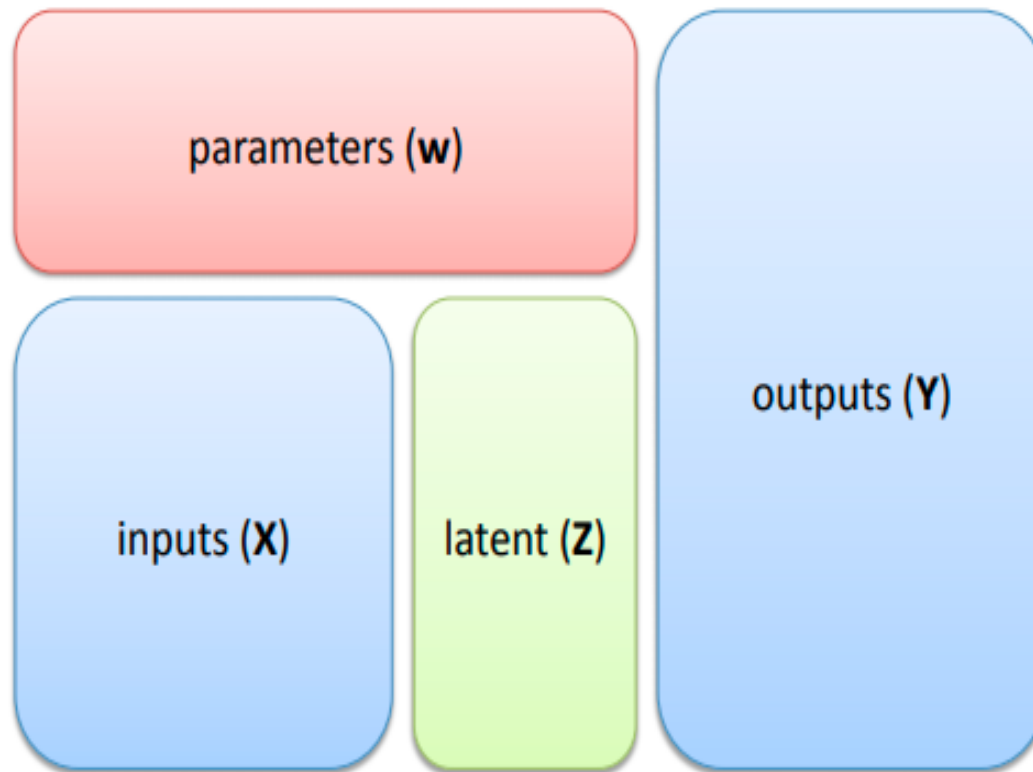
**Swati Pratap Jagdale**

Department of Computer Science and Engineering

### Hidden variables in Decoding and Supervised Learning



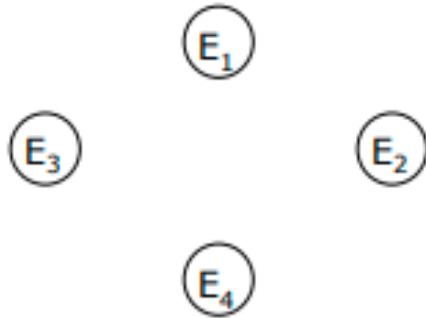
### Hidden variables in Decoding and Supervised Learning



## Learning With Hidden Variables

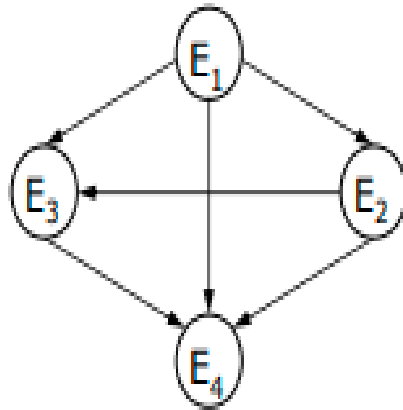
---

- Consider a situation in which you can observe a whole bunch of different evidence variables,  $E_1$  through  $E_n$ . Maybe they're all the different symptoms that a patient might have. Or maybe they represent different movies and whether someone likes them.



## Learning With Hidden Variables

- If those variables are all conditionally dependent on one another, then we'd need a highly connected graph that's capable of representing the entire joint
- distribution between the variables. Because the last node has  $n-1$  parents, it will take on the order of  $2^n$  parameters to specify the conditional probability tables in this network.



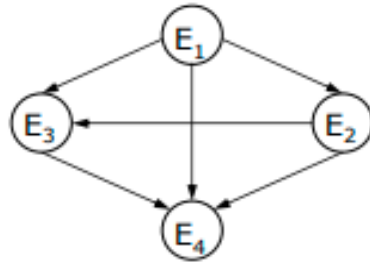
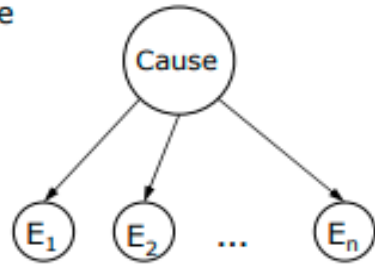
$O(2^n)$  parameters

Without the cause,  
all the evidence is  
dependent on  
each other

## Learning With Hidden Variables

- But, in some cases, we can get a considerably simpler model by introducing an
- additional “cause” node. It might represent the underlying disease state that was causing the patients’ symptoms or some division of people into those who like westerns and those who like comedies.

Cause is unobservable



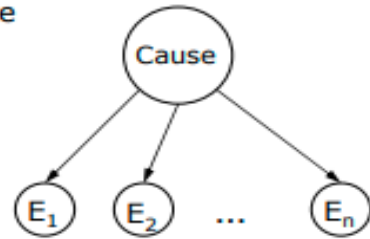
Without the cause,  
all the evidence is  
dependent on  
each other

$O(2^n)$  parameters

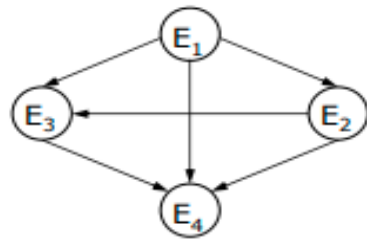
## Learning With Hidden Variables

- In the simpler model, the evidence variables are conditionally independent given the causes. That means that it would only require on the order of  $n$  parameters to describe all the CPTs in the network, because at each node, we just need a table of size 2 (if the cause is binary; or  $k$  if the cause can take on  $k$  values), and one (or  $k-1$ ) parameter to specify the probability of the cause.

Cause is unobservable



$O(n)$  parameters



$O(2^n)$  parameters

Without the cause,  
all the evidence is  
dependent on  
each other



## Learning With Hidden Variables

---

- In the simpler model, the evidence variables are conditionally independent given the causes. That means that it would only require on the order of  $n$  parameters to describe all the CPTs in the network, because at each node, we just need a table of size 2 (if the cause is binary; or  $k$  if the cause can take on  $k$  values), and one (or  $k-1$ ) parameter to specify the probability of the cause.
- So, what if you think there's a hidden cause? How can you learn a network with unobservable variables?

## Simpson's Paradox

---

- Edward Hugh Simpson, a statistician and former cryptanalyst at Bletchley Park, described the statistical phenomenon - Simpson's paradox
- The art of data science is seeing beyond the data — using and developing methods and tools to get an idea of what that hidden reality looks like.
- Simpson's paradox showcases the importance of skepticism and interpreting data with respect to the real world, and also the dangers of oversimplifying a more complex truth by trying to see the whole story from a single data-viewpoint.

## Simpson's Paradox

---

- ***Simpson's Paradox:***

***A trend or result that is present when data is put into groups that reverses or disappears when the data is combined.***

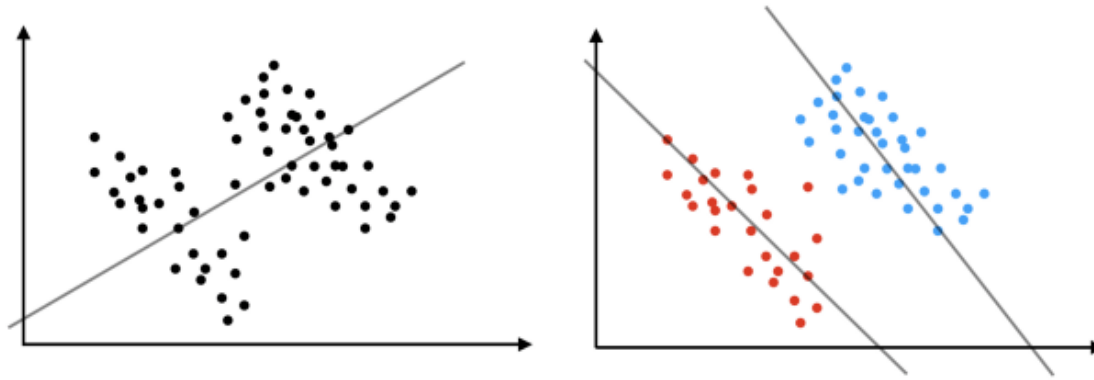
Example:

UC Berkley's suspected gender-bias.

At the beginning of the academic year in 1973, UC Berkeley's graduate school had admitted roughly 44% of their male applicants and 35% of their female applicants.

- There was a statistically significant gender bias in favour of women for 4 out of the 6 departments, and no significant gender bias in the remaining 2.
- It is discovered that women tended to apply to departments that admitted a smaller percentage of applicants overall, and that this hidden variable affected the marginal values for the percentage of accepted applicants in such a way as to reverse the trend that existed in the data as a whole.

## Simpson's Paradox



A visual example: the overall trend reverses when data is grouped by some colour-represented category.

## Simpson's Paradox

a simple example in business:

- Suppose we're in the soft drinks industry and we're trying to choose between two new flavours we've produced. We could sample public opinion on the two flavours

Flavour	Sample Size	# Liked Flavour
Sinful Strawberry	1000	800
Passionate Peach	1000	750

- 80% of people enjoyed 'Sinful Strawberry' whereas only 75% of people enjoyed 'Passionate Peach'. So 'Sinful Strawberry' is more likely to be the preferred flavour.

## Simpson's Paradox

- Some other information while conducting the survey, such as the sex of the person sampling the drink. What happens if we split our data up by sex?
- 84.4% of men and 40% of women liked 'Sinful Strawberry' whereas 85.7% of men and 50% of women liked 'Passionate Peach'

Flavour	# Men	# Liked Flavour (Men)	# Women	# Liked Flavour (Women)
Sinful Strawberry	900	760	100	40
Passionate Peach	700	600	300	150

## Simpson's Paradox

- *according to our sample data, generally people prefer 'Sinful Strawberry', but both men and women separately prefer 'Passionate Peach'.*
- *This is an example of Simpson's Paradox!*

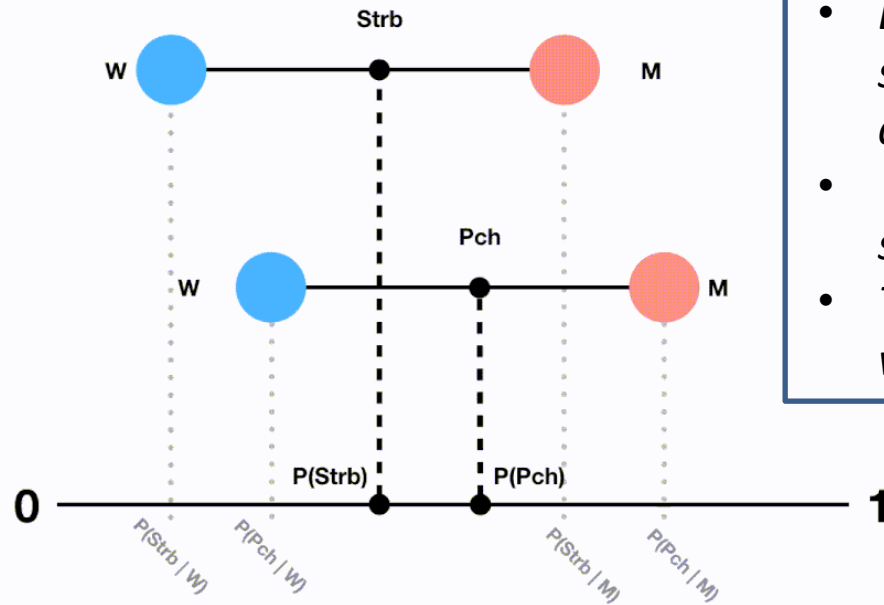
Flavour	# Men	# Liked Flavour (Men)	# Women	# Liked Flavour (Women)
Sinful Strawberry	900	760	100	40
Passionate Peach	700	600	300	150

- **Lurking variables (Hidden Variables)**
- Simpson's paradox arises when there are hidden variables that split data into multiple separate distributions.
- Such a hidden variable is aptly referred to as a **lurking variable**, and they can often be difficult to identify.
- **Consider the lurking variable (sex) and a little bit of probability theory:-**
- $P(\text{Liked Strawberry}) = P(\text{Liked Strawberry} \mid \text{Man})P(\text{Man}) + P(\text{Liked Strawberry} \mid \text{Woman})P(\text{Woman})$
- $800/1000 = (760/900) \times (900/1000) + (40/100) \times (100/1000)$
- $P(\text{Liked Peach}) = P(\text{Liked Peach} \mid \text{Man})P(\text{Man}) + P(\text{Liked Peach} \mid \text{Woman})P(\text{Woman})$
- $750/1000 = (600/700) \times (700/1000) + (150/300) \times (300/1000)$



## Simpson's Paradox

- **Lurking variables (Hidden Variables)**
- We can think of the marginal probabilities of sex ( $P(\text{Man})$  and  $P(\text{Woman})$ ) as weights that, in the case of 'Sinful Strawberry', cause the total probability to be significantly shifted towards the male opinion.



- Each coloured circle represents either the men or women that sampled each flavour, the position of the centre of each circle corresponds to that group's probability of liking the flavour.
- As the circles grow (i.e. sample proportions change) we can see how the marginal probability of liking the flavour changes.
- The marginal distributions shift and switch as samples become weighted with respect to the lurking variable (sex).

## References

---

<http://www.cs.cmu.edu/~nasmith/psnlp/lecture5.pdf>

<https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-825-techniques-in-artificial-intelligence-sma-5504-fall-2002/lecture-notes/Lecture18FinalPart1.pdf>

<https://towardsdatascience.com/simpsons-paradox-and-interpreting-data-6a0443516765>



**THANK YOU**

---

**Swati Pratap Jagdale**

Department of Computer Science

[swatigambhire@pes.edu](mailto:swatigambhire@pes.edu)