

ISE 540

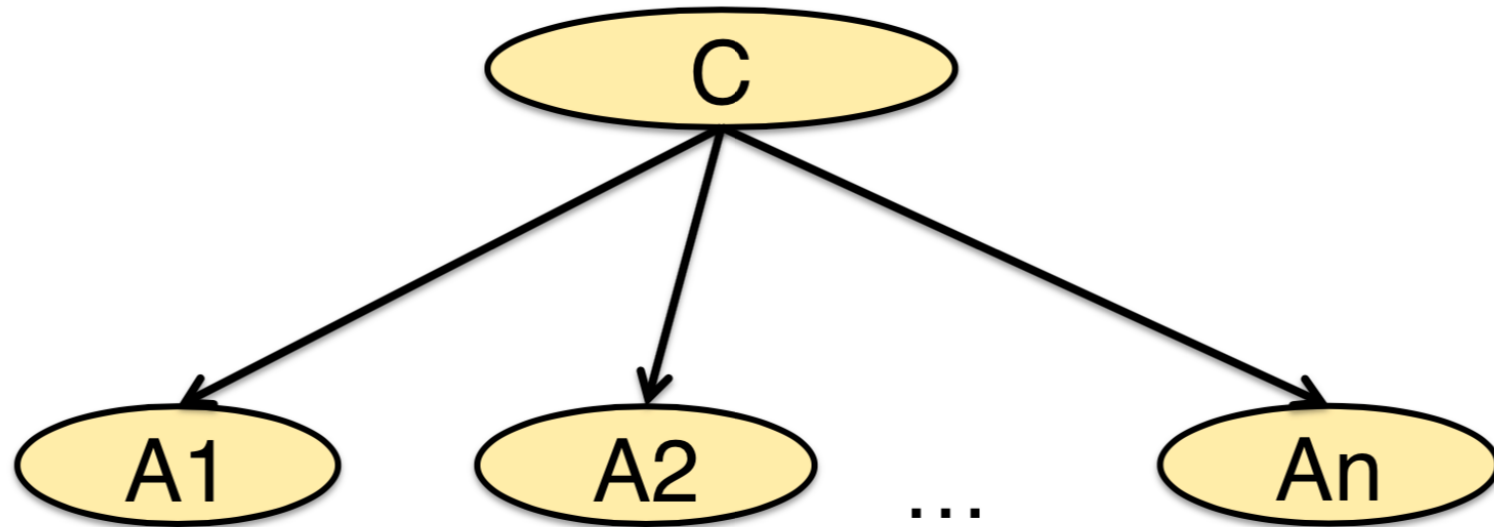
Text Analytics

Mayank Kejriwal

Research Assistant Professor/Research Lead
Department of Industrial and Systems Engineering
Information Sciences Institute
USC Viterbi School of Engineering

kejriwal@isi.edu

Naïve Bayes (NB)



Each item has a number of attributes

$$A_1=a_1, \dots, A_n=a_n$$

We predict the class c based on

$$c = \operatorname{argmax}_c \prod_i P(A_i = a_i \mid C=c) P(C=c)$$

Does the customer want sugar?

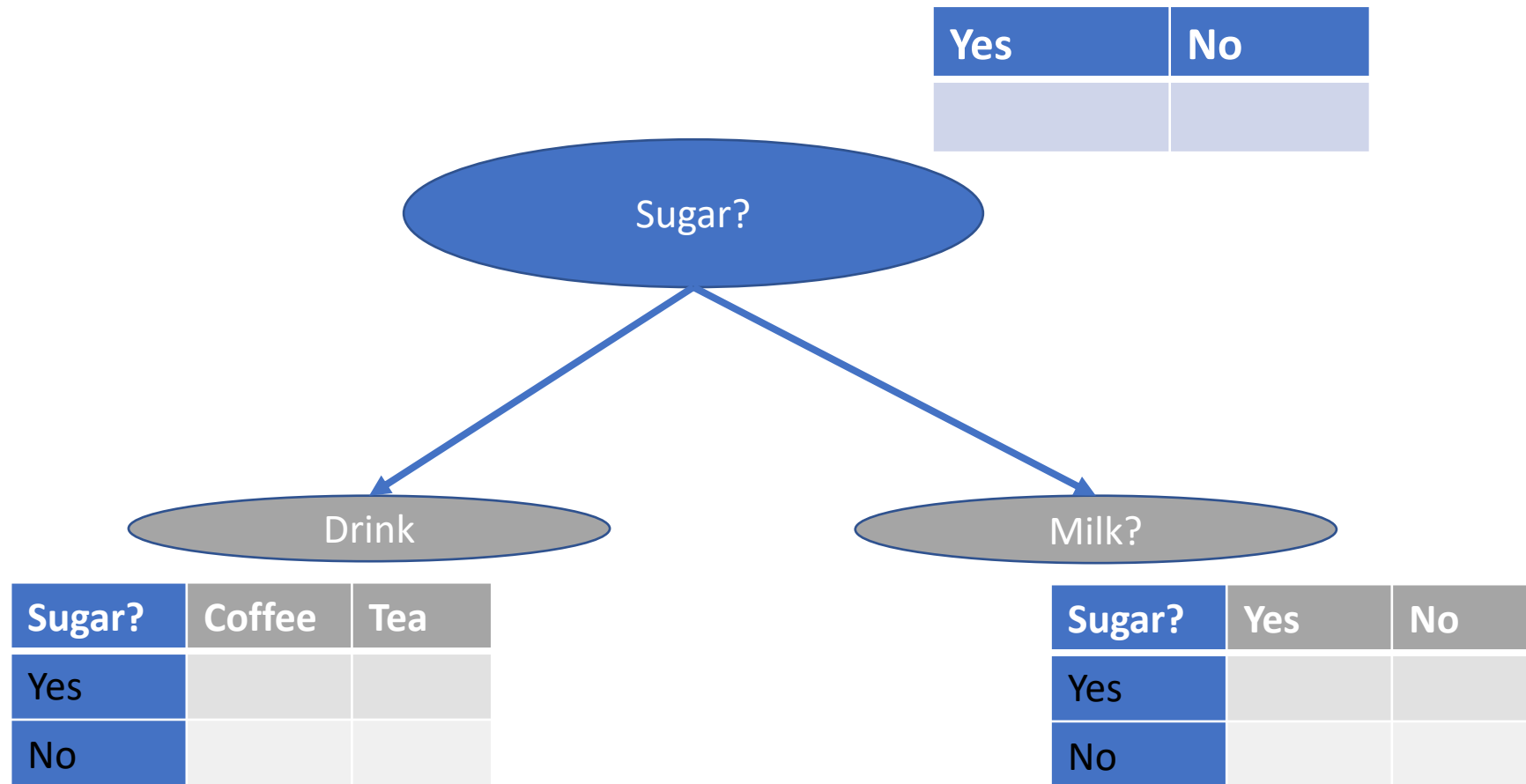
x1	x2	Y
A1: drink	A2: milk?	C: sugar?
coffee	no	yes
coffee	yes	no
tea	yes	yes
tea	no	no

Can you train a Naïve Bayes classifier to predict whether the customer wants sugar or not?

What is $P(\text{coffee} \mid \text{sugar})$?

Getting the NB parameters

- Use maximum likelihood!



Does the customer want sugar?

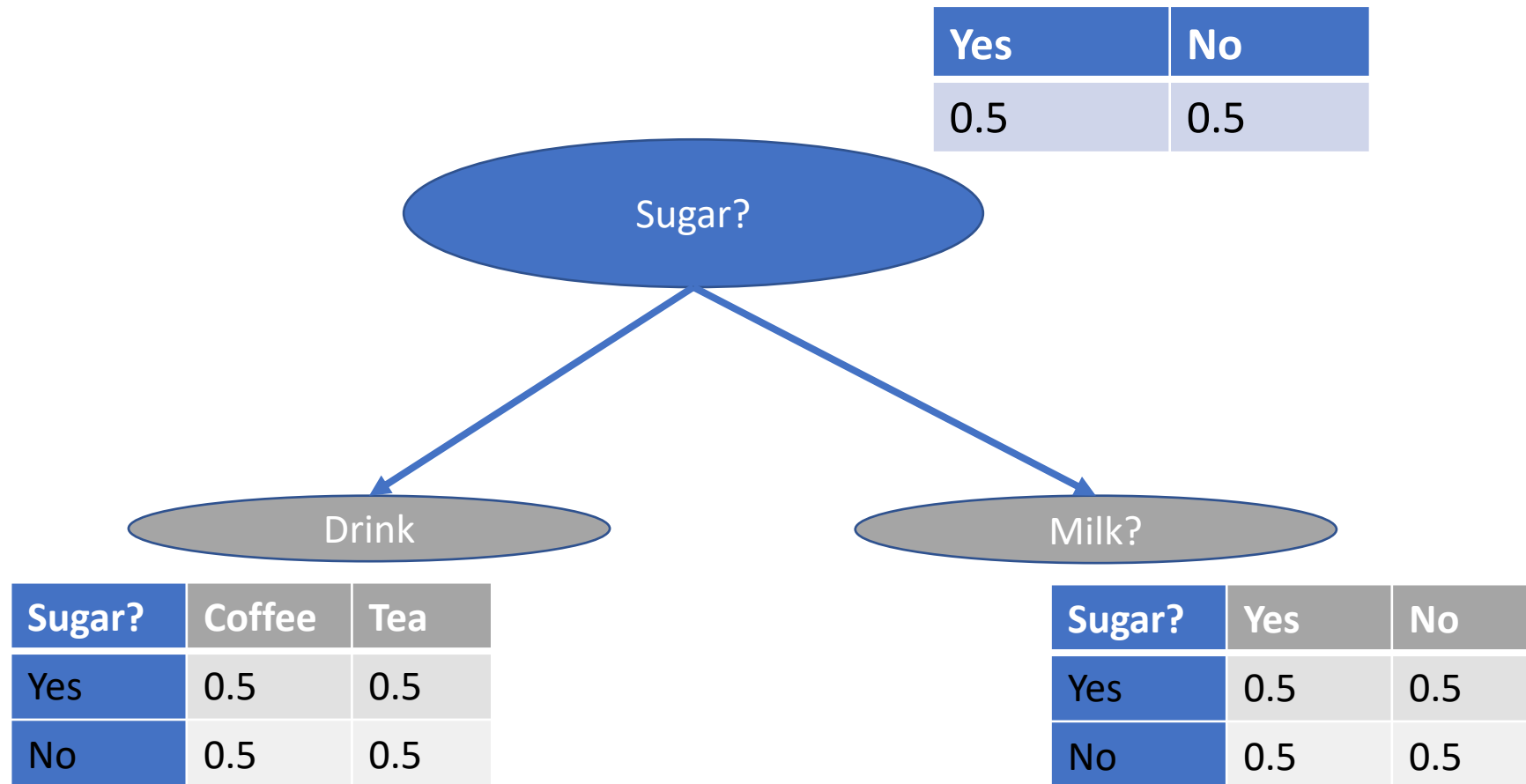
x1	x2	Y
A1: drink	A2: milk?	C: sugar?
coffee	no	yes
coffee	yes	no
tea	yes	yes
tea	no	no

Can you train a Naïve Bayes classifier to predict whether the customer wants sugar or not?

What is $P(\text{coffee} \mid \text{sugar})$?

Getting the NB parameters

- Use maximum likelihood!



Another example

x1	x2	Y
A1: drink	A2: milk?	C: sugar?
coffee	no	yes
coffee	yes no	no
tea	yes	yes
tea coffee	no	no

Can you train a Naïve Bayes classifier to predict whether the customer wants sugar or not?

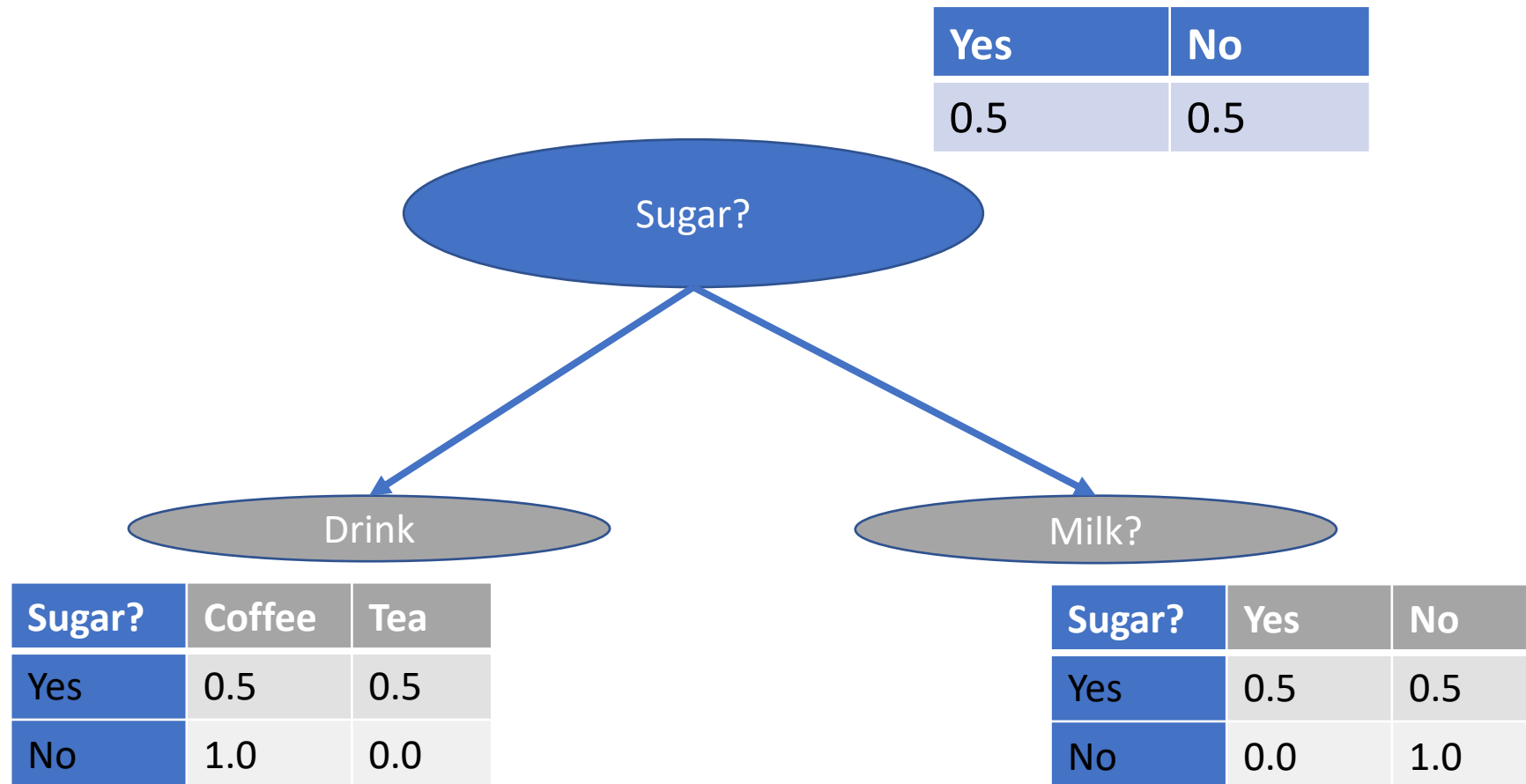
What is $P(\text{coffee} \mid \text{sugar})$?

Why is this a problem?

Conditional Independence: $P(x_1, x_2 \mid y) = P(x_1 \mid y) P(x_2 \mid y)$

For example: $P(\text{drink}=\text{coffee}, \text{milk}=\text{yes} \mid \text{sugar}=\text{yes}) = P(\text{drink}=\text{coffee} \mid \text{sugar}=\text{yes})$

$P(\text{milk}=\text{yes} \mid \text{sugar}=\text{yes})$



Prediction problem

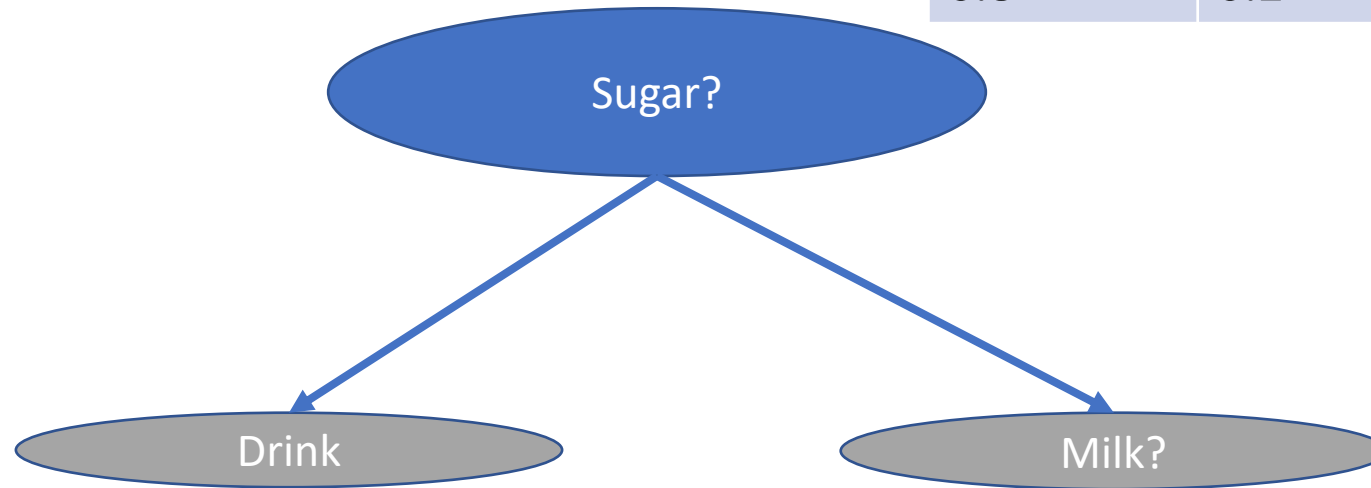
- Given milk and drink, can you predict whether I want sugar or not?
- We want to pick the value for sugar (yes or no) that maximizes the probability $P(\text{sugar} \mid \text{milk}, \text{drink})$
- From standard probability theory:
 - $P(\text{sugar} \mid \text{milk}, \text{drink}) = P(\text{milk}, \text{drink} \mid \text{sugar})P(\text{sugar})$
- Using conditional independence assumption:
 - $P(\text{milk}, \text{drink} \mid \text{sugar}) = P(\text{milk} \mid \text{sugar}) P(\text{drink} \mid \text{sugar})$
- Putting the two together,
 - $P(\text{sugar} \mid \text{milk}, \text{drink}) = P(\text{milk} \mid \text{sugar}) P(\text{drink} \mid \text{sugar}) P(\text{sugar})$
- But these are just the 'cells' in the NB table!
- What we want to do is to compute the 'argmax' for this expression over the 'sugar' variable

Another example

--What does NB predict given that I'm having coffee with milk? What is the probability of the prediction?

--What about tea without milk?

Yes	No
0.8	0.2



Sugar?	Coffee	Tea
Yes	0.6	0.4
No	0.9	0.1

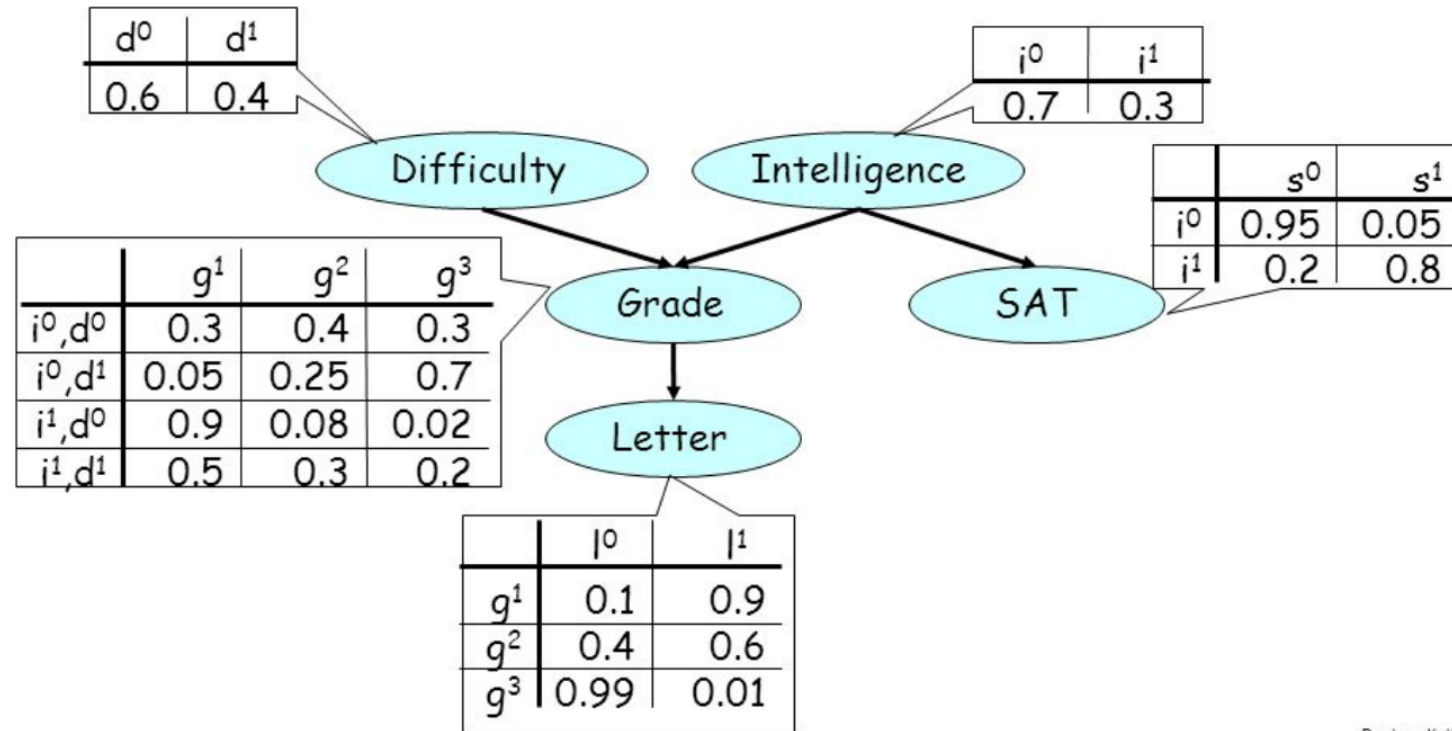
Sugar?	Yes	No
Yes	0.3	0.7
No	0.6	0.4

What can go wrong ?

Think conditional independence assumption...

Can we make this model more
'general'?

The Student Network



Daphne Koller

Difficulty (of the class): Takes values 0 (low difficulty) and 1 (high difficulty)

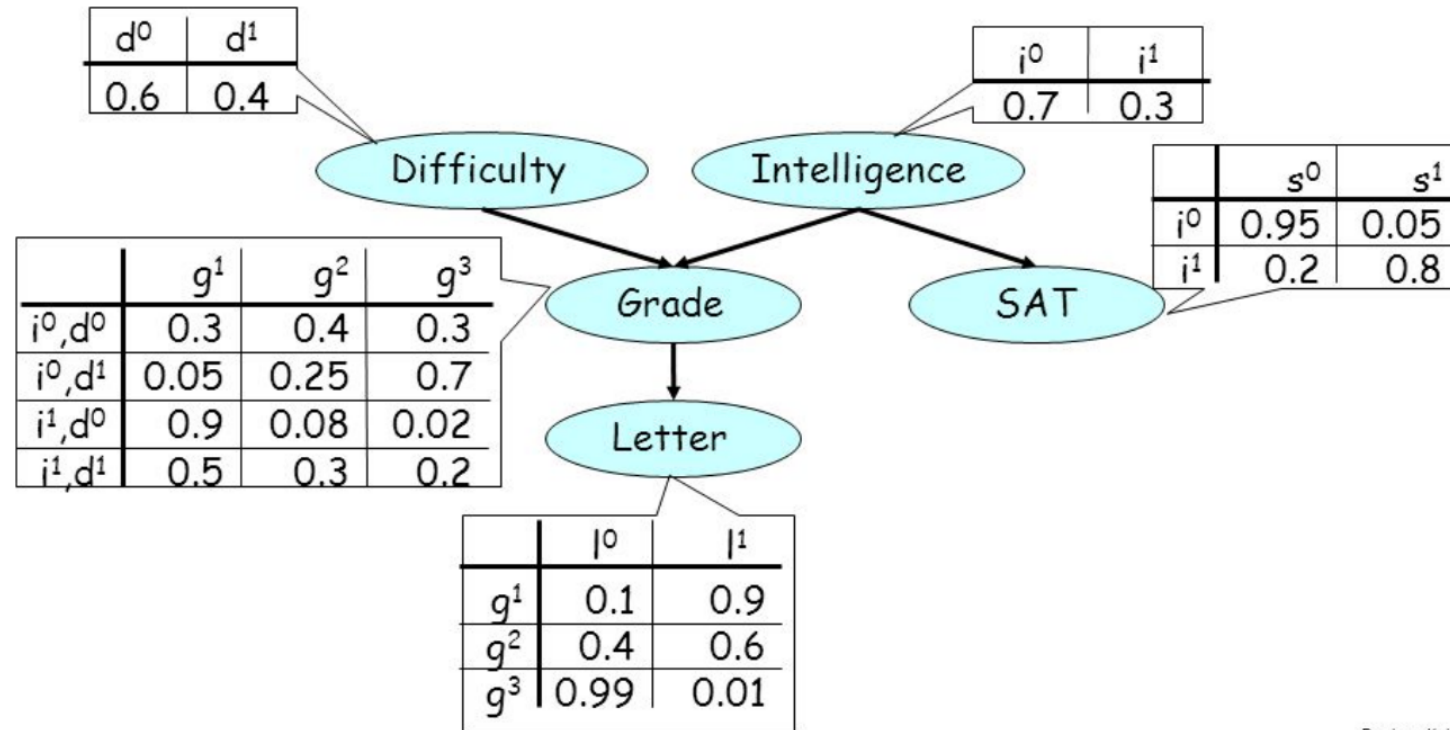
Intelligence (of the student): Takes values 0 (not intelligent) and 1 (intelligent)

Grade (the student gets in the class): Takes values 1 (good grade), 2 (average grade), and 3 (bad grade)

SAT (student's score in the SAT exam): Takes values 0 (low score) and 1 (high score)

Letter (quality of recommendation letter the student gets from the professor after completing the course): Takes values 0 (not a good letter) and 1 (a good letter)

The Student Network



Daphne Koller

What is the probability that a student gets a good letter (l^1) given that the student's grade is g^2 ?

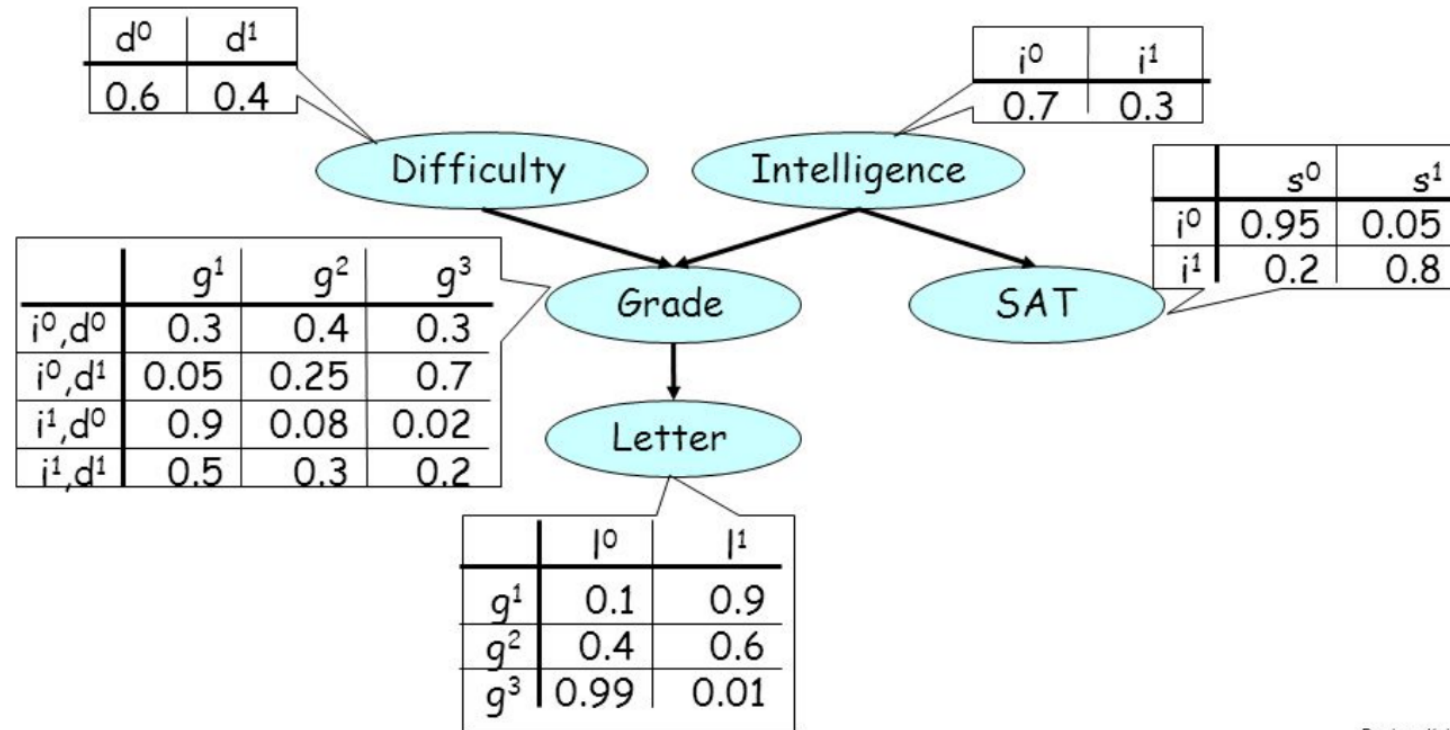
Difficulty (of the student): Takes values 0 (not intelligent) and 1 (intelligent)

Grade (the student gets in the class): Takes values 1 (good grade), 2 (average grade), and 3 (bad grade)

SAT (student's score in the SAT exam): Takes values 0 (low score) and 1 (high score)

Letter (quality of recommendation letter the student gets from the professor after completing the course): Takes values 0 (not a good letter) and 1 (a good letter)

The Student Network



Daphne Koller

Does it matter whether the student is 'intelligent' or not?

NO iff $P(l^1 | i^1) = P(l^1 | i^0)$, YES iff $P(l^1 | i^1) \neq P(l^1 | i^0)$

Difficulty (of the task): Takes values 0 (not difficult) and 1 (difficult)

Intelligence (of the student): Takes values 0 (not intelligent) and 1 (intelligent)

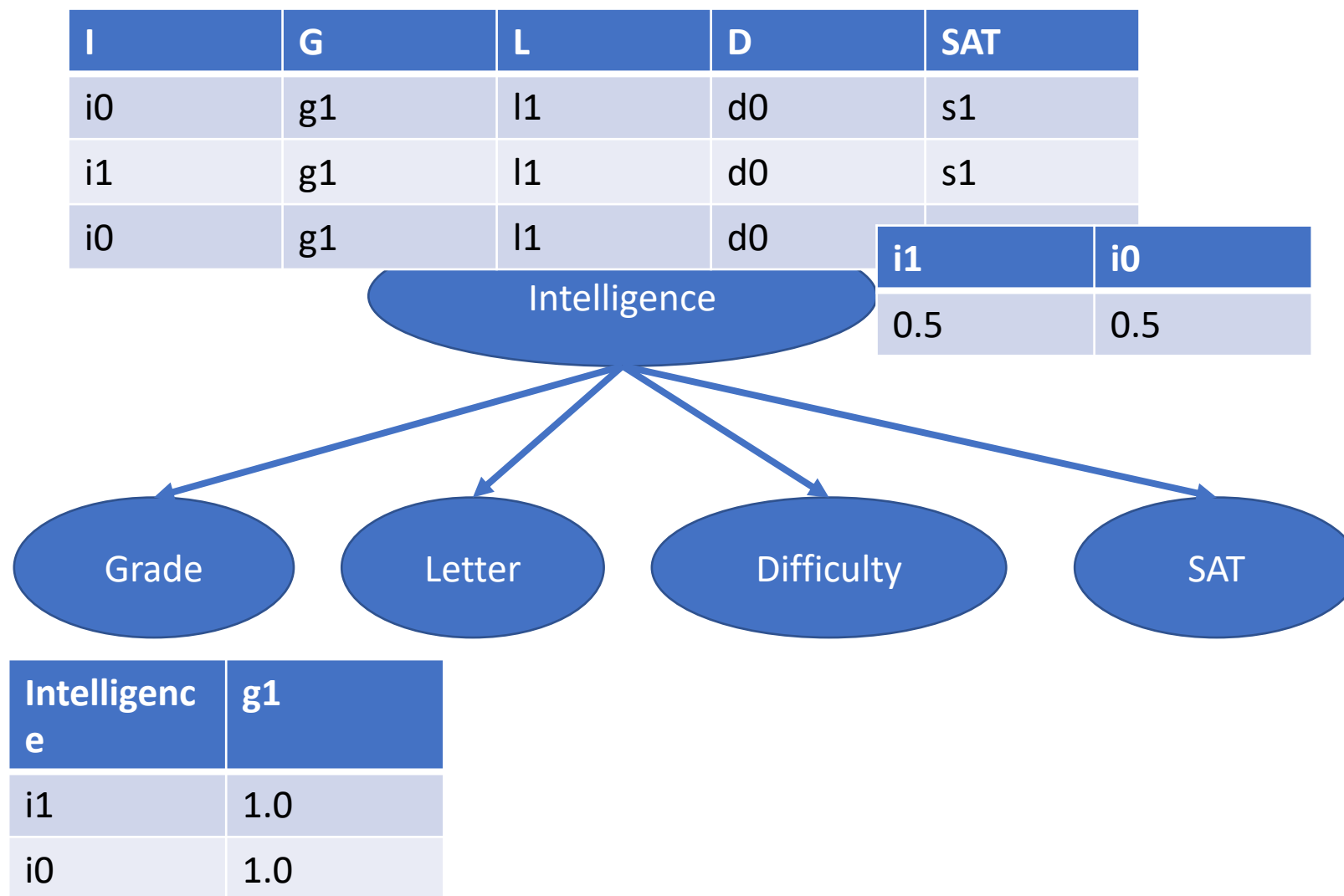
Grade (the student gets in the class): Takes values 1 (good grade), 2 (average grade), and 3 (bad grade)

SAT (student's score in the SAT exam): Takes values 0 (low score) and 1 (high score)

Letter (quality of recommendation letter the student gets from the professor after completing the course): Takes values 0 (not a good letter) and 1 (a good letter)

How do we represent the student's network as a naïve Bayes?

- First question you must ask is what is the 'class' variable (usually the variable you will be trying to 'predict')...
 - Suppose it is Intelligence



What if the target class was 'Letter'? What would be the NB model? Can the two models give 'different' results?

If I gave you a 'table of observations' just like the coffee/sugar example, would you be able to infer the CPTs?