

# GATE in Data Science & Artificial Intelligence

## ARTIFICIAL INTELLIGENCE

Introduction to reasoning  
under uncertainty

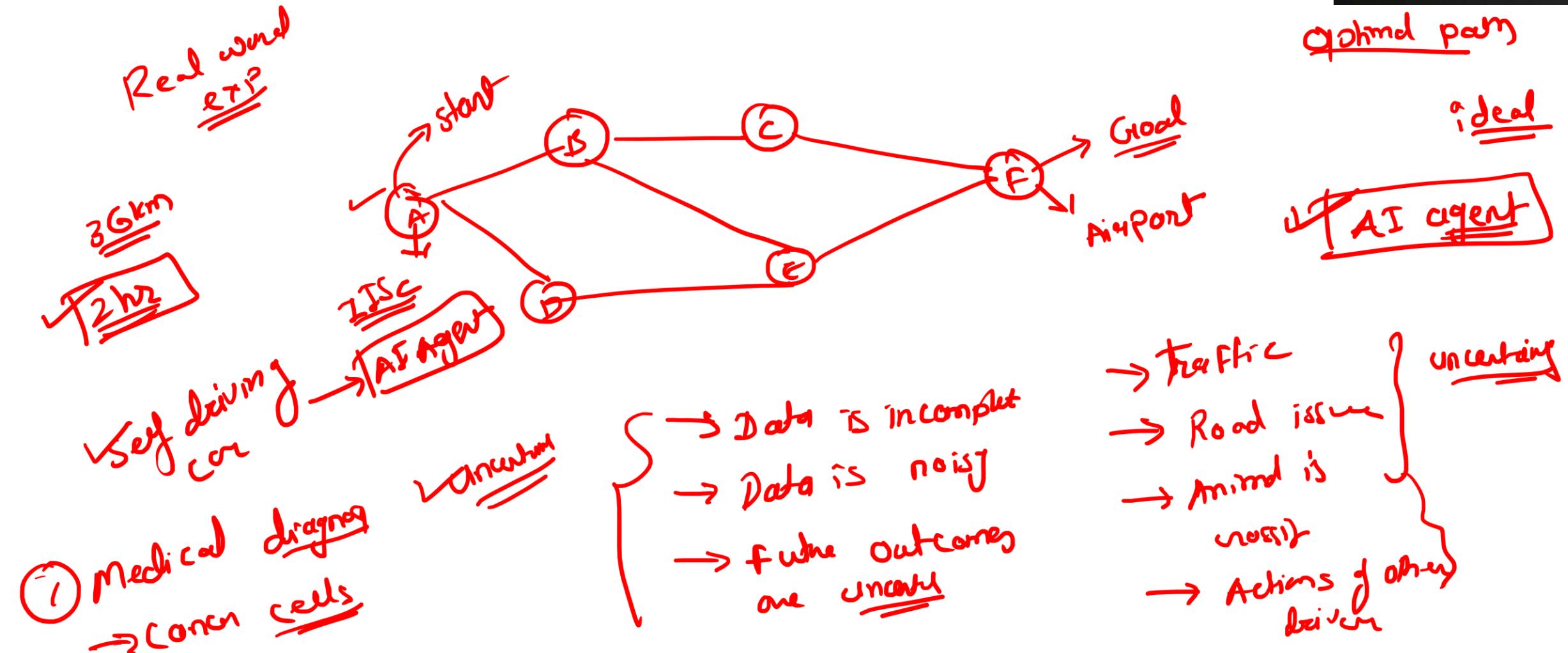
By:

Piyush Wairale

MTech (IIT Madras)

Instructor at IIT Madras BS in Data Science Degree





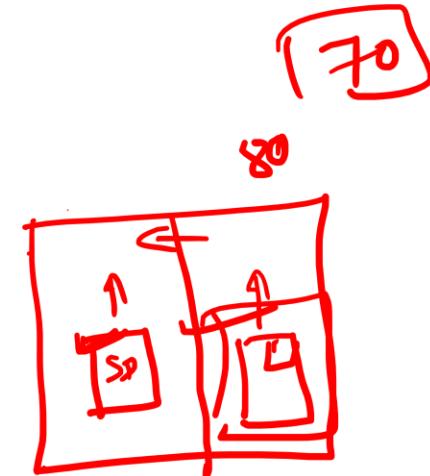


Robot

Reasoning under uncertainty is a fundamental concept in artificial intelligence (AI) that involves making decisions or drawing conclusions when the information available is incomplete, ambiguous, or subject to change. Unlike deterministic systems, where outcomes are predictable and certain, AI systems often operate in environments where uncertainty is a significant factor.

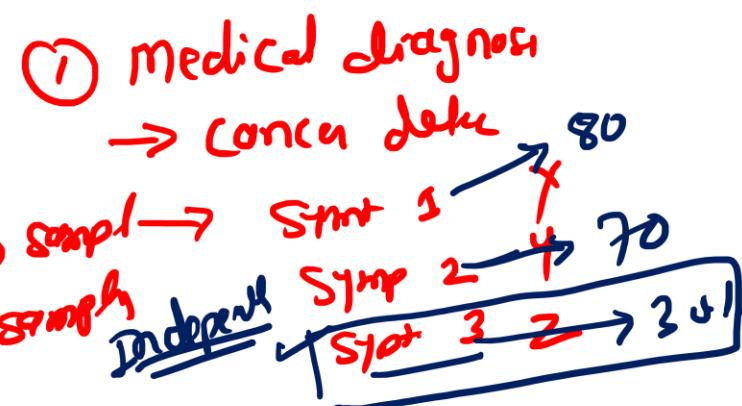
① Probabilistic Models → Bayes Meets  
 → Bayesian Network → graphical Models  
 → Markov Random Field  
 → Conditional Independence of Probabilistic Model

$$P(A|B)$$



② Inference :-  
 → Exact Inf :- involve computing exact prob  
 → using algm of variable elimination

② Approximate Inf :-  
 → when exact inf is computationally infeasible, then  
 we will go for approximate Inf  
 → using Sampling → Monte Carlo sample → Gibbs sample



# Bayesian Network



Piyush Wairale

⇒ Bayes theorem (conditional probability)  
(Belief Network)

A graphical probabilistic Models  
that represent a set of variable  
and their conditional dependencies.  
using DAG

✓ A (Disease)

B (Test)

✓ T

+ve

✓ T

+ve

✓ T

-ve

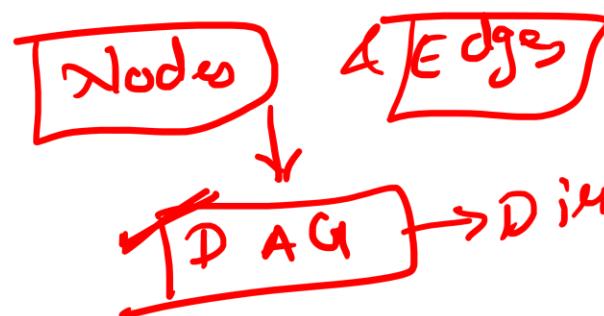
F

-ve

✓ F

+ve

✓ F



Conditioned Independent  
Bayesian Network

Nodes :- variable → CRV  
Edges :- conditional dependencies → DRV

Prior  $P(A) = 0.5$   
 $P(A = \text{diseas}) = 0.5$   
 $P(A = \text{not diseas}) = 1 - 0.5 = 0.5$

$\sqrt{B}$  is conditional dependent on A

$P(B|A)$



Piyush Wairale

$P(B = +ve | A = T)$   
 $P(B = -ve | A = T)$   
 $P(B = +ve | A = F)$   
 $P(B = -ve | A = F)$



## 1. Nodes (Variables):

- Each node in a Bayesian Network represents a random variable, which could be discrete or continuous. The variable can represent anything from observable quantities (e.g., sensor readings) to latent variables (e.g., underlying causes that are not directly observed).

## 2. Edges (Dependencies):

- Directed edges (arrows) between nodes represent conditional dependencies between the variables. If there is a directed edge from node A to node B, it means that B is conditionally dependent on A.



## 3. Conditional Probability Distributions (CPDs):

- Associated with each node is a conditional probability distribution that quantifies the effect of the parent nodes (those with directed edges leading to the node) on the node itself. For a node  $X_i$  with parents  $Parents(X_i)$ , the CPD specifies  $P(X_i | Parents(X_i))$ .

$$P(B|A)$$



## 4 Joint Probability Distribution:

- The joint probability distribution over all the variables in the network can be expressed as the product of the conditional probabilities for each variable given its parents:

$$\checkmark P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(\underline{X_i} \mid \overset{\text{Parents}}{\uparrow} Parents(X_i))$$

D (Disease) :- True / False

T (Test) :- {Positive / Negative}

S (Symptom) :- (Present / Absent)

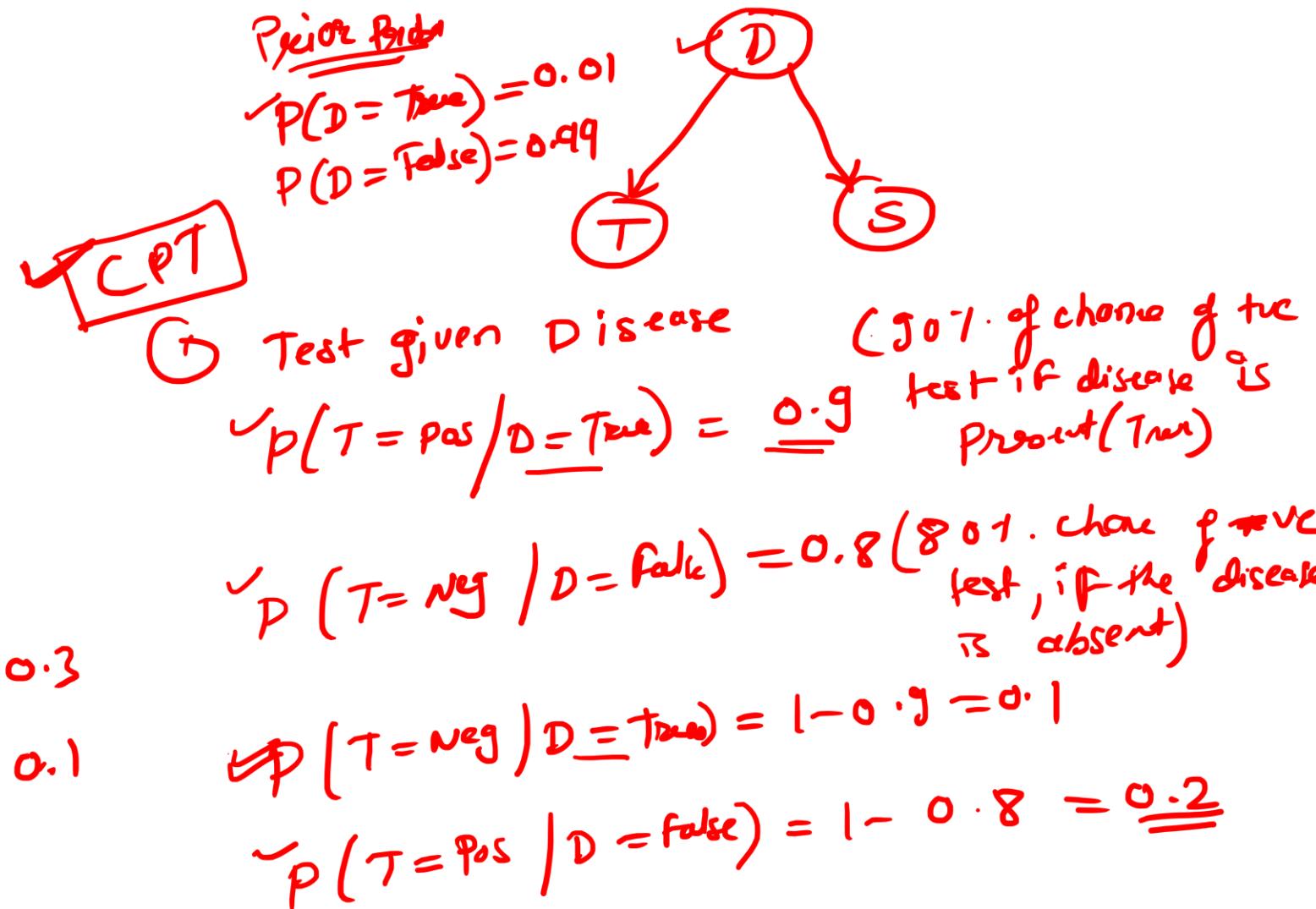
② Symptom given Disease

$$P(S = \text{Present} | D = \text{True}) = 0.7$$

$$P(S = \text{Absent} | D = \text{True}) = 0.3$$

$$\checkmark P(S = \text{Absent} | D = \text{False}) = 1 - 0.7 = 0.3$$

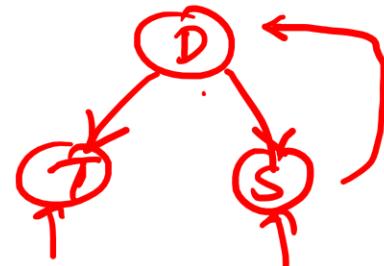
$$\checkmark P(S = \text{Present} | D = \text{False}) = 1 - 0.3 = 0.1$$





\* Joint Prob Dist

$$P(D, T, S) = P(D) \times P(T|D) \times P(S|D)$$



$D=?$ ,  $T=+ve$ ,  $S=Present$

$\Rightarrow$  the prob that the person having a disease  $\rightarrow D = True$

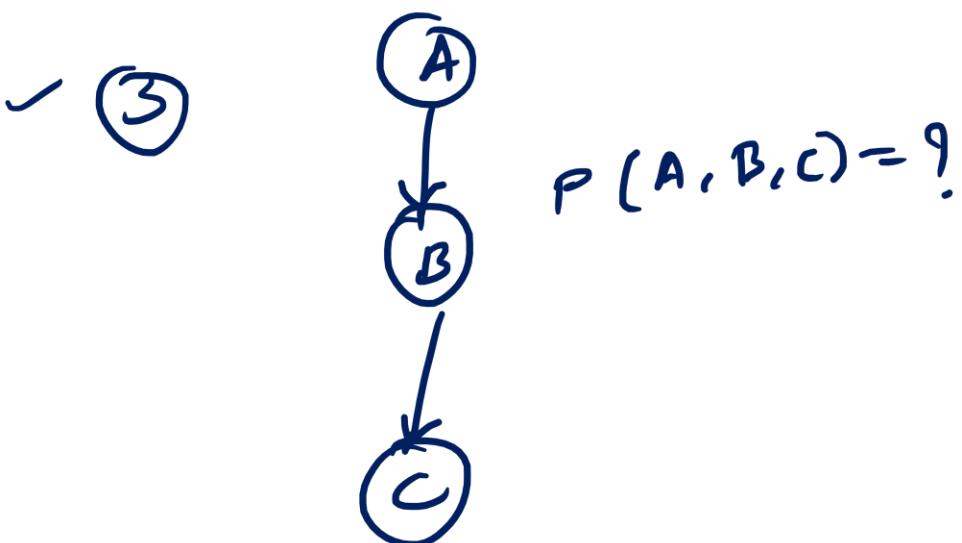
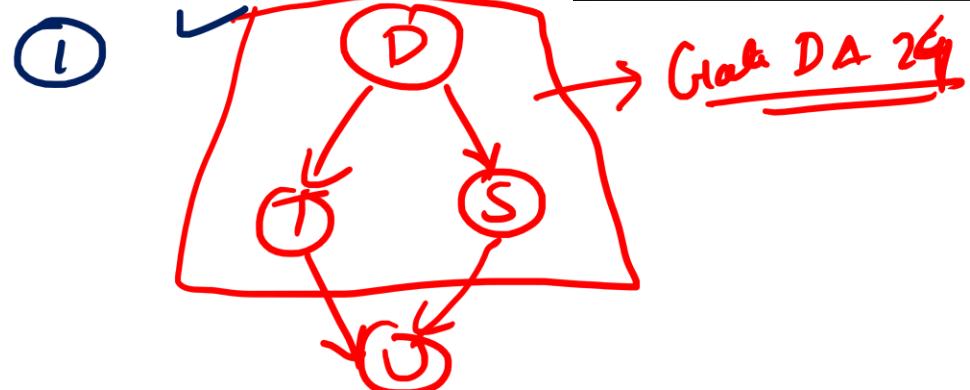
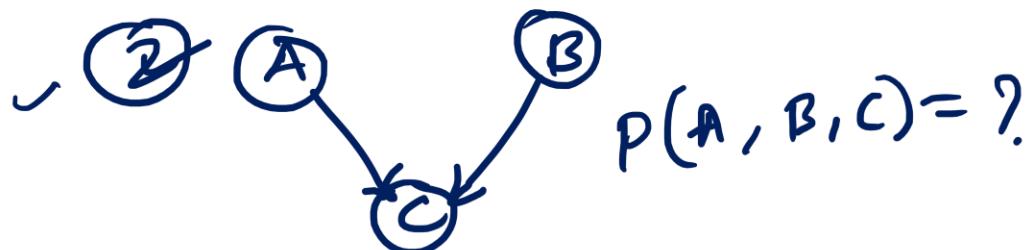
$$\Rightarrow P(D=True, T=+ve, S=True) = P(D=True) \times P(T=+ve|D=True) \times P(S=True|D=True)$$

$$\begin{aligned} P(D=True, T=+ve, S=True) &= 0.01 \times 0.9 \times 0.7 \\ &= \underline{\underline{0.0063}} \end{aligned}$$

$$\begin{aligned} \textcircled{3} P(D=True, T=-ve, S=False) &= ? \end{aligned}$$

$$\textcircled{2} P(D=True, T=-ve, S=False) = ?$$

# Bayesian Network



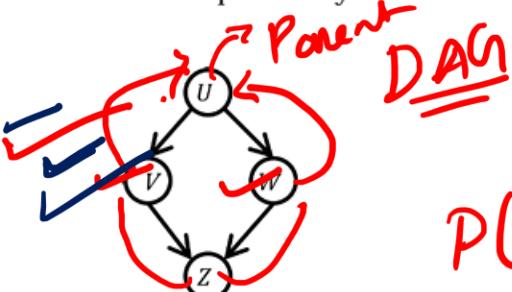


Q.64

AI

Given the following Bayesian Network consisting of four Bernoulli random variables and the associated conditional probability tables:

2m



$$P(V|U)$$

$$P(W|U)$$

CPT

	$P(\cdot)$
$U = 0$	0.5
$U = 1$	0.5

	$P(V = 0   \cdot)$	$P(V = 1   \cdot)$
$U = 0$	0.5	0.5
$U = 1$	0.5	0.5
	$P(W = 0   \cdot)$	$P(W = 1   \cdot)$
$U = 0$	1	0
$U = 1$	0	1

		$P(Z = 0   \cdot)$	$P(Z = 1   \cdot)$
$V = 0$	$W = 0$	0.5	0.5
$V = 0$	$W = 1$	1	0
$V = 1$	$W = 0$	1	0
$V = 1$	$W = 1$	0.5	0.5

The value of  $P(U = 1, V = 1, W = 1, Z = 1) = \underline{\hspace{2cm}}$  (rounded off to three decimal places).

Joint Probability :- for given network

$$P(U, V, W, Z) = P(U) \times P(V|U) \times P(W|U) \times P(Z|V, W)$$

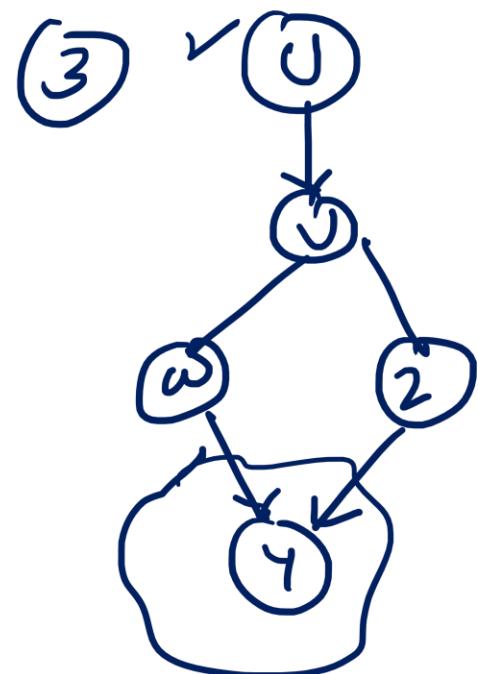
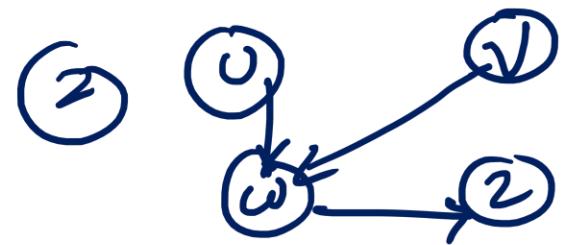
$$P(U=1, V=1, W=1, Z=1) = P(U=1) \times P(V=1|U=1) \times P(W=1|U=1) \times P(Z=1|V=1, W=1)$$

$$= 0.5 \times 0.5 \times 1 \times 0.5$$

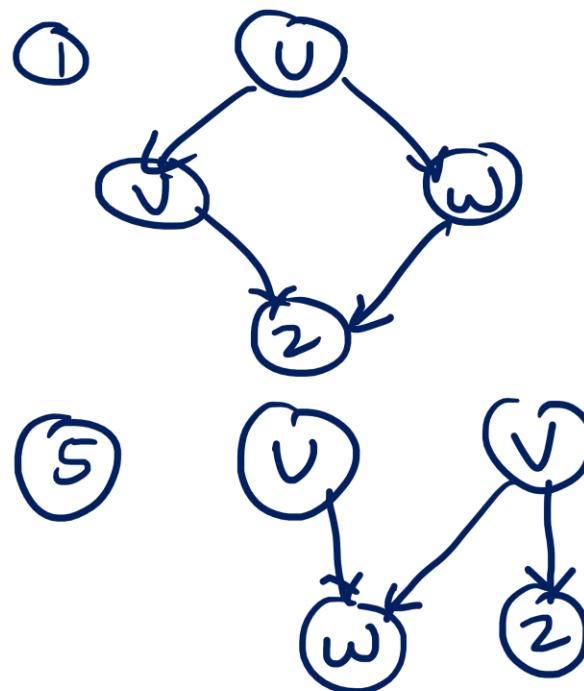
$$= 0.125$$

$$(1) P(U=1, V=0, W=0, Z=1)$$

$$(2) P(U=0, V=1, W=1, Z=1)$$



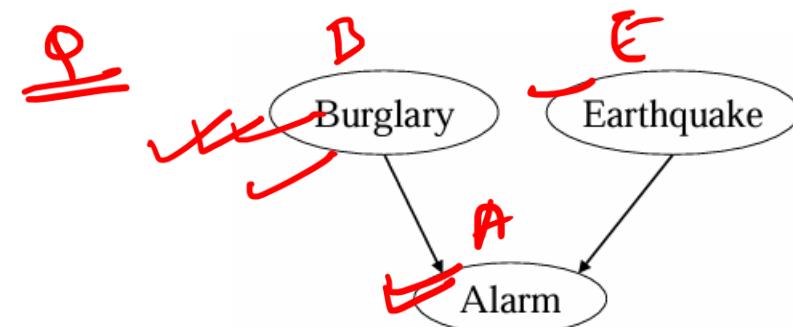
④



③ Nodes  
⑤ Edges

- ① DAG
- ② CPT
- ③ Required condit. Prob

# Bayesian Network Example



B	P(B)	E	P(E)
false	0.999	false	0.998
true	0.001	true	0.002

$$\begin{aligned}
 P(B, E, A) &= P(B) \cdot P(E) \cdot P(A|B, E) \\
 &= P(B=f) \cdot P(E=T) \cdot P(A=T|B=f, E=T) \\
 &= 0.999 \times 0.002 \times 0.29 \\
 &= ?
 \end{aligned}$$

$\left\{ \begin{array}{l} B=\text{false} \\ E=\text{True} \\ A = \underline{\text{True}} \end{array} \right.$

B	E	A	P(A B,E)
false	false	false	0.999
false	false	true	0.001
false	true	false	0.71
false	true	true	0.29
true	false	false	0.06
true	false	true	0.94
true	true	false	0.05
true	true	true	0.95

- HwE
- (1)  $B = \text{True}, E = \text{false}, A = \text{True}$
  - (2)  $B = \text{F}, E = \text{T}, A = \text{True}$
  - (3)  $B = \text{F}, E = \text{T}, A = \text{False}$

# Bayesian Network Example



Q

P(B)	
true	false
0.001	0.999

2

Burglary

2

Earthquake

$P(J|A, B, E)$

Alarm

$2^3 = 8$

B	E	$P(A B, E)$	
		true	false
true	true	0.95	0.05
true	false	0.94	0.06
false	true	0.29	0.71
false	false	0.001	0.999

A	$P(J A)$	
	true	false
true	0.9	0.1
false	0.05	0.95

John calls

$2^2 = 4$

Mary calls

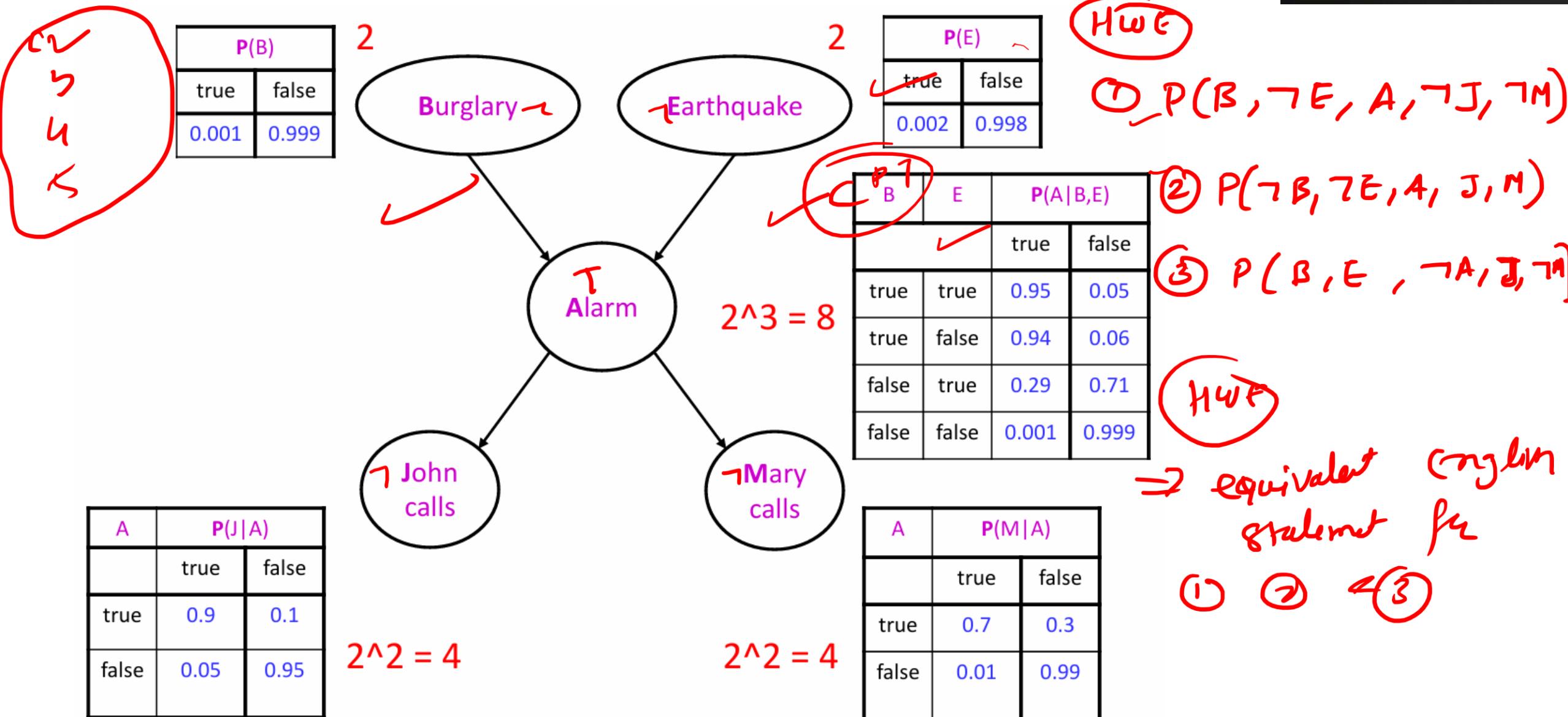
$2^2 = 4$

A	$P(M A)$	
	true	false
true	0.7	0.3
false	0.01	0.99

$$\begin{aligned}
 & \frac{P(B, E, A, J, M)}{P(B) \cdot P(E)} \times P(A|BE) \\
 & \times P(J|A) \times P(M|A) \\
 & = 0.001 \times 0.002 \\
 & \times 0.95 \times 0.9 \\
 & \times 0.7 \\
 & = 
 \end{aligned}$$

$$\begin{aligned}
 P(A, B, E) &= P(B) \cdot P(E) \\
 &\times P(A|B, E) \times P(J|A, B, E) \\
 &\times P(M|A, B, E)
 \end{aligned}$$

# Bayesian Network Example



# Bayesian Network Example

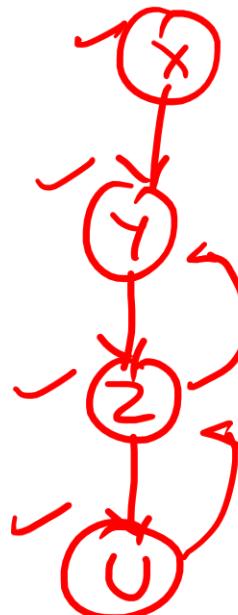


$$P(x, y, z, u)$$

$$= P(x) \cdot P(y|x) \cdot P(z|y, x) \\ \times P(u|z, y, x)$$

by using chain rule

$$\checkmark = P(x) \cdot P(y|x) \cdot P(z|y) \\ P(u|z)$$



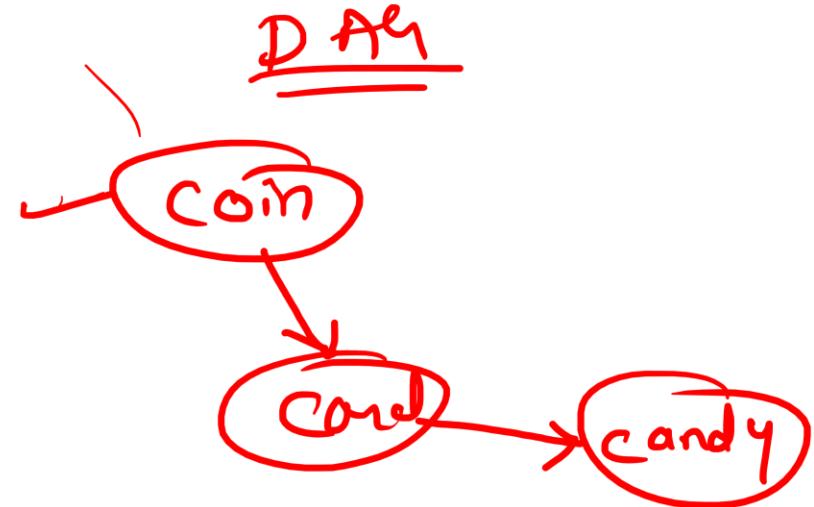
# Bayesian Network Example



Coin	P(Coin)
tails	0.5
heads	0.5

Coin	Card	P(Card   Coin)
tails	black	0.6
tails	red	0.4
heads	black	0.3
heads	red	0.7

Card	Candy	P(Candy   Card)
black	1	0.5
black	2	0.2
black	3	0.3
red	1	0.1
red	2	0.3
red	3	0.6

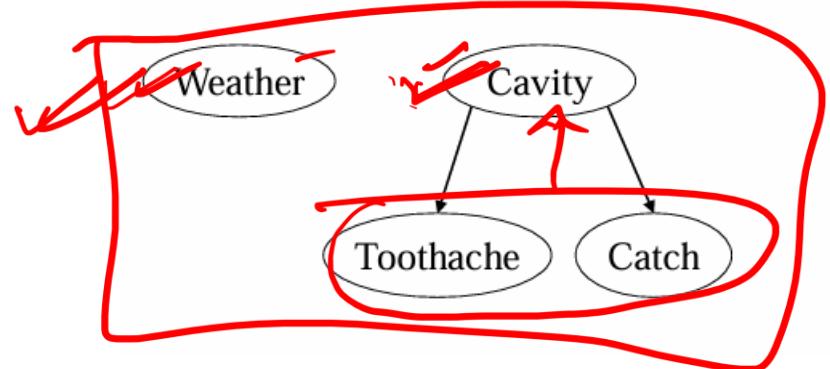


What does the DAG for this Bayes net look like?

$$P(c, \text{card}, \text{candy}) = P(c) \cdot (P(\text{card}|c)) \\ \times P(\text{candy}|\text{card})$$



## Bayesian Network Example



$$P(w, c, T, \text{catch}) = P(w) \cdot P(c) \cdot P(T|c) \\ + P(\text{catch}|\text{cavity})$$

$$= P(c, T|\text{catch}) \times P(w) \boxed{P(\text{cavity}|T, \text{catch})}$$

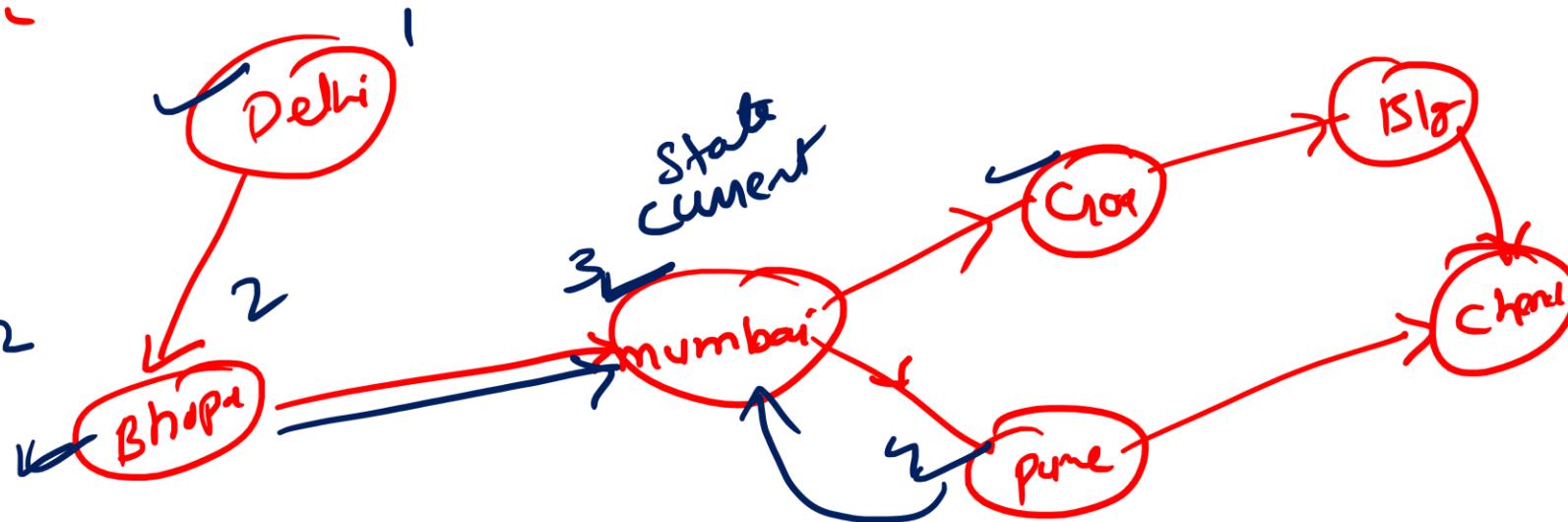
① weather & cavity are absolute independent

② Toothache and catch are conditional independent given cavity



\* Bayesian Network

① Today  
 $P(\text{Rain}) = 0.8$   
 $P(\text{Sunny}) = 0.2$



$p(\text{Rain} = \text{true})$

$P(m|B, D)$

$x_4 = \underline{\underline{x_3}}$

$$P(D, B, m, P) = P(D) \cdot P(B|D) \times P(m|B, D) \\ \times P(pure|m, B, D)$$

$$\boxed{P(P|m)}$$



A **Markov Chain** is a mathematical system that undergoes transitions from one state to another within a finite or countable number of possible states. It is a type of stochastic process that satisfies the Markov property, which states that the future state depends only on the current state and not on the sequence of events that preceded it.

## Key Components:

1. **States:** The set of possible conditions in which the system can exist.
2. **Transition Probabilities:** The probabilities of moving from one state to another.
3. **Initial State Distribution:** The probability distribution over the states at the start of the process. *current*

## Markov Property

The basic property of a Markov chain is that only the most recent point in the trajectory affects what happens next.

p<sub>nm</sub>

This is called the Markov Property.

It means that  $X_{t+1}$  depends upon  $X_t$ , but it does not depend upon  $X_{t-1}, \dots, X_1, X_0$ .

$x_{n+1}$   
 $x_n$   
 $x_{n-1}$   
 $x_{n-2}$   
 $x_0$

1m



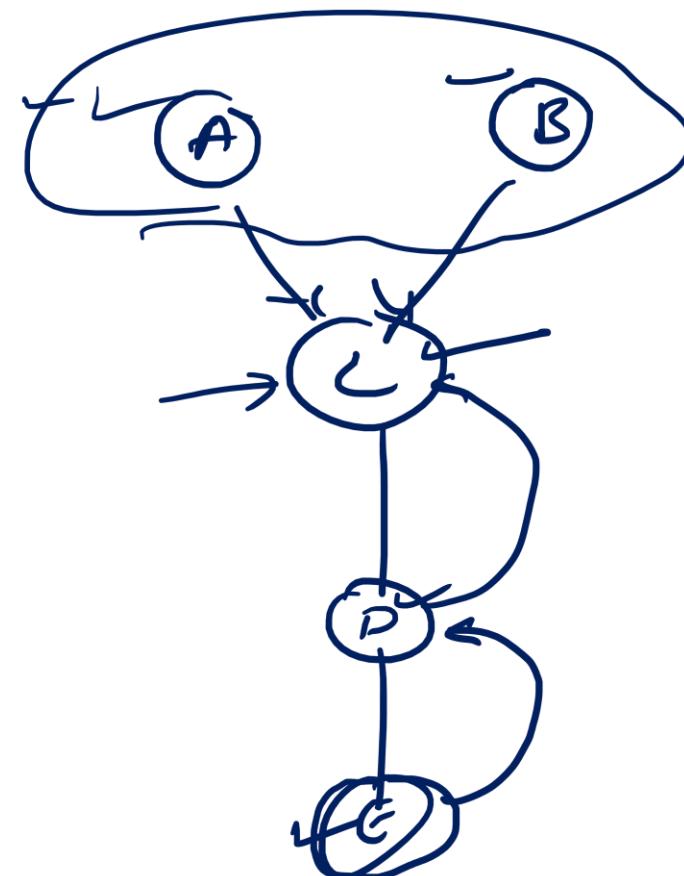
We formulate the Markov Property in mathematical notation as follows:

$$\mathbb{P}(X_{t+1} = s | X_t = s_t, X_{t-1} = s_{t-1}, \dots, X_0 = s_0) = \mathbb{P}(X_{t+1} = s | X_t = s_t),$$

for all  $t = 1, 2, 3, \dots$  and for all states  $s_0, s_1, \dots, s_t, s$ .

*Explanation:*

$\mathbb{P}(X_{t+1} = s | X_t = s_t, X_{t-1} = s_{t-1}, X_{t-2} = s_{t-2}, \dots, X_1 = s_1, X_0 = s_0)$   
 distribution of  $X_{t+1}$       depends on  $X_t$   
 but whatever happened before time  $t$   
 doesn't matter.





Example Weather forecast

Today:- Day 1  $\rightarrow$  comes start

If it's sunny today, there is chance of 80% that it will be sunny tomorrow  
 If it's rainy today, there is chance of 60% that it will rain tomorrow

T States Sunny (S), Rainy (R)

Tomorrow  $\rightarrow$  Transition Probabilities :-

$$\textcircled{1} \quad P(S_{tom} | S_{today}) = 0.8$$

$$\textcircled{2} \quad P(R_{tom} | S_{today}) = 0.2$$

$$\textcircled{3} \quad P(S_{tom} | R_{today}) = 0.4$$

$$\textcircled{4} \quad P(R_{tom} | R_{today}) = 0.6$$

Day 2

Day 4

That it will be sunny tomorrow

that it will rain tomorrow

Day after tomorrow

If it's Sunny tomorrow

$$P(Sunny = DFT | Sunny_{tom}) = 0.8 \times 0.8 \\ = 0.64$$

If it's Rainy tomorrow

$$P(Sunny = DFT | R_{tom}) = 0.4 \times 0.2 \\ = 0.08$$

Fwd Prob of Sunny 2 day from now

$$P(Sunny \text{ in 2 days}) = 0.64 + 0.08 \\ = 0.72$$



A **Hidden Markov Model (HMM)** is an extension of a **Markov Chain** where the **states** are not **directly observable** (hidden). Instead, each state produces an **observable output** (or emission) according to a probability distribution.

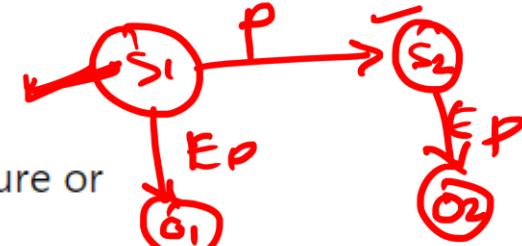
## Key Components:

1. **States**: The hidden states of the system (e.g., underlying weather conditions like high pressure or low pressure).

2. **Transition Probabilities**: The probabilities of transitioning from one hidden state to another.

3. **Emission Probabilities**: The probabilities of observing a certain output given the current hidden state.

4. **Initial State Distribution**: The probability distribution over the hidden states at the start.



# Hidden Markov Model (HMM)

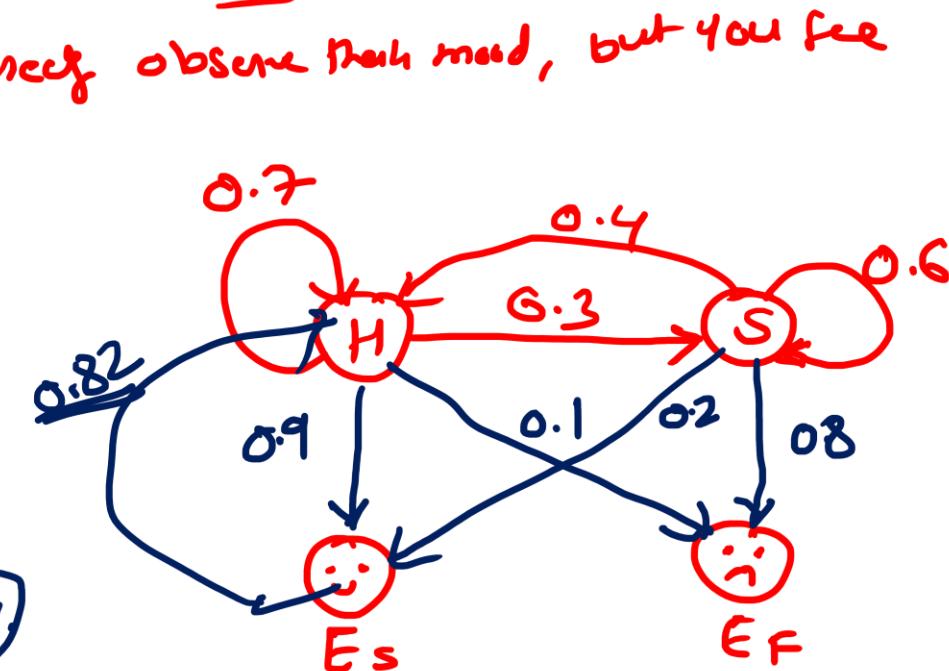


\* We want to guess your bf/gf mood based on their behavior.  
The mood can be Happy or Sad, but you can't directly observe their mood, but you see their actions like Smiling or Frowning.

⇒ Hidden states    ① Happy (H)    ② Sad (S)

⇒ Observation :- Smiling (Smile) 😊  
Frowning (frown) 😕

Q You observe that your gf/bf is Smiling today, what can you infer about their mood



$$\Rightarrow P(\text{smile} | H) = 0.9$$

$$(P(\text{smile} | S)) = 0.2$$



\* Initial Prob. :  $P(H) = 0.5$        $P(S) = 0.5$

$$P(\text{😊} | H) = \frac{P(S \cap H)}{P(H)}$$

$\neg P(H \cap \text{😊}) = 0.5 \times 0.9 = 0.45 =$

$\vee P(S \cap \text{😊}) = 0.5 \times 0.2 = 0.1$

\* Normalize Prob.

① Total prob that you are bf/gf & smily =  $0.45 + 0.1 = 0.55$

$$P(H|\text{😊}) = \frac{0.45}{0.55} = 0.82 \checkmark$$

$$P(S|\text{😊}) = \frac{0.1}{0.55} = 0.182$$



Hidden Markov Models (HMMs) are a type of statistical model used to represent systems that are assumed to be a Markov process with unobservable (hidden) states. HMMs are particularly useful in scenarios where you have observable data that is believed to be influenced by hidden states, such as in speech recognition, natural language processing, and bioinformatics.

Three fundamental algorithms associated with HMMs are:

1. Forward Algorithm
2. Viterbi Algorithm
3. Baum-Welch Algorithm



## 1. Forward Algorithm

The Forward Algorithm is used to compute the **probability of a sequence of observations** given a Hidden Markov Model. It helps determine how likely a particular sequence of observed events is, given the model parameters.

*Decoding*

## 2. Viterbi Algorithm

The Viterbi Algorithm is used to find ~~the most likely sequence of hidden states~~ (called the Viterbi path) that could have generated a given sequence of observations. This is useful in applications like speech recognition, where you want to determine the most probable sequence of spoken words given an audio signal.



### 3. Baum-Welch Algorithm

- The **Baum-Welch Algorithm** is a special case of the **Expectation-Maximization (EM) algorithm** used to find the unknown parameters of an HMM. It is used for **training HMMs** when the state transition probabilities and emission probabilities are not known.
- The Baum-Welch Algorithm, also known as the Expectation-Maximization (EM) algorithm for HMMs, is used to **learn the parameters** (transition probabilities, emission probabilities, and initial state probabilities) of an HMM from a set of observed sequences when the states are unknown. It finds the **parameters** that maximize the likelihood of the observed sequences.



## Summary of Each Algorithm

- **Forward Algorithm:** Computes the probability of an observation sequence given the HMM parameters.
- **Viterbi Algorithm:** Finds the most likely sequence of hidden states for a given observation sequence.
- **Baum-Welch Algorithm:** Estimates the HMM parameters from observed sequences using an iterative approach.



## ~~Inference~~ in HMMs:

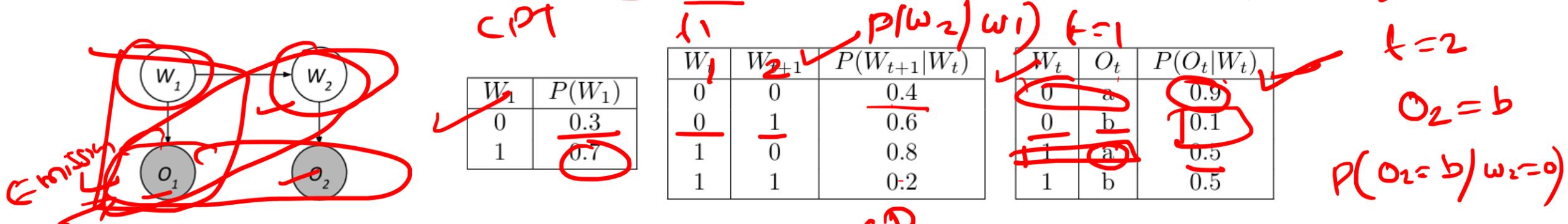
The two main tasks are:

1. ~~Decoding~~: Given a sequence of observations, determine the most likely sequence of hidden states. This is often done using the ~~Viterbi algorithm~~.
2. ~~Learning~~: Estimate the model parameters (transition and emission probabilities) given a set of observed sequences. This is typically done using the ~~Baum-Welch algorithm~~.

# Hidden Markov Model (HMM)



Consider the following Hidden Markov Model.  $O_1$  and  $O_2$  are supposed to be shaded.



Suppose that we observe  $O_1 = a$  and  $O_2 = b$ .

Using the forward algorithm, compute the probability distribution  $P(W_2 | O_1 = a, O_2 = b)$  one step at a time.

**JP**

$$P(w_2, o_1 = a, o_2 = b) = P(o_2 | w_2) \cdot P(w_2 | o_1)$$

$$\equiv P(o_2 = b | w_2) \cdot P(w_2 | o_1 = a)$$

$$P(w_1, o_1 = a) = P(w_1) \cdot P(o_1 | w_1)$$

$$\Rightarrow P(w_1 = 0, o_1 = a) = P(w_1 = 0) \cdot P(o_1 = a | w_1 = 0)$$

$$= 0.3 \times 0.9 = 0.27$$

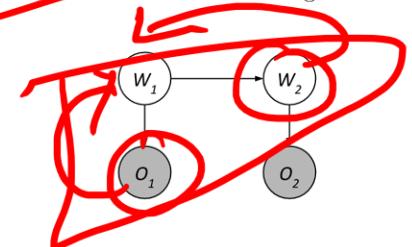
$$\Rightarrow P(w_1 = 1, o_1 = a) = P(w_1 = 1) \cdot P(o_1 = a | w_1 = 1)$$

$$= 0.7 \times 0.5 = 0.35$$

# Hidden Markov Model (HMM)



✓ Consider the following Hidden Markov Model.  $O_1$  and  $O_2$  are supposed to be shaded.



$W_1$	$P(W_1)$
0	0.3
1	0.7

$W_t$	$W_{t+1}$	$P(W_{t+1} W_t)$
0	0	0.4
0	1	0.6
1	0	0.8
1	1	0.2

$W_t$	$O_t$	$P(O_t W_t)$
0	a	0.9
0	b	0.1
1	a	0.5
1	b	0.5

$$\begin{aligned} & P(w_2, O_1=a) \\ & P(w_1=0, O_1=a) = \\ & P(w_2=1, O_1=a) = \end{aligned}$$

Suppose that we observe  $O_1 = a$  and  $O_2 = b$ .

Using the forward algorithm, compute the probability distribution  $P(W_2|O_1 = a, O_2 = b)$  one step at a time.

$$\begin{aligned} P(w_2=0|O_1=a) &= \sum_{w_1=0,1} P(w_1, O_1=a) \cdot P(w_2|w_1) \xrightarrow{\text{total pm}} \\ &= P(w_1=0, O_1=a) \cdot P(w_2=0|w_1=0) + P(w_1=1, O_1=a) \cdot P(w_2=0|w_1=1) \\ &= 0.27 \times 0.4 + 0.35 \times 0.8 = 0.388 \end{aligned}$$

$$\begin{aligned} P(w_2=1|O_1=a) &= \sum_{w_1=0,1} P(w_1, O_1=a) \cdot P(w_2=1|w_1) \xrightarrow{\text{total pm}} \\ &= P(w_1=0, O_1=a) \cdot P(w_2=1|w_1=0) + P(w_1=1, O_1=a) \cdot P(w_2=1|w_1=1) \\ &= 0.27 \times 0.6 + 0.35 \times 0.2 = 0.232 \end{aligned}$$



$$P(\omega_2, o_1=a, o_2=b) = P(\omega_2, o_1=a) \cdot P(o_2=b | \omega_2)$$

①  $\omega_2=0$

$$\begin{aligned} P(\omega_2=0, o_1=a, o_2=b) &= P(\omega_2=0, o_1=a) \cdot P(o_2=b | \omega_2=0) \\ &= 0.388 \times 0.1 = 0.0388 \end{aligned}$$

②  $\omega_2=1$

$$\begin{aligned} P(\omega_2=1, o_1=a, o_2=b) &= \frac{P(\omega_2=1, o_1=a)}{P(o_2=b | \omega_2=1)} \\ &= 0.232 \times 0.5 = 0.116 \end{aligned}$$

by Normalization

$$\Rightarrow P(\omega_2 | o_1=a, o_2=b)$$

$$\textcircled{1} \quad P(\omega_2=0 | o_1=a, o_2=b) = \frac{0.0388}{0.0388 + 0.116} = 0.25$$



$$\Rightarrow P(\omega_2=1 \mid o_1=a, o_2=b) = \underline{P(\omega_2=1, o_1=a, o_2=b)}$$

$$= \frac{0.116}{0.388 + 0.116} = \underline{\underline{0.75}}$$



Probabilistic inference involves computing probabilities and making predictions based on a probabilistic model. It is widely used in fields such as statistics, machine learning, artificial intelligence, and data science to draw conclusions from data that may have uncertainty or noise.

Probabilistic inference can be broadly divided into two types: exact inference and approximate inference.

$$P(X=t|e) = ?$$



Computing  $P(X = x|e)$  in a GM is an NP-hard problem, which means that there is not generalized efficient algorithm for inference given an arbitrary PGM, query and evidence nodes. However, for some families of models, there are polynomial time algorithms for inference. Another solution to deal with the hardness of the problem is finding approximate solutions. The following is a list of some of the exact and approximate algorithms on graphical models.



Exact inference refers to the process of computing the exact probability distribution or the exact values of interest from a probabilistic model. This involves finding the precise solution to a problem, considering all possible scenarios or configurations of the random variables.

## Examples of Exact Inference:

using variable elimination

• Bayesian Networks: Exact inference in Bayesian networks involves computing posterior probabilities of certain variables given evidence. Algorithms like Variable Elimination and the Junction Tree Algorithm are commonly used for this purpose.

• Markov Chains: For discrete-time Markov chains, exact inference involves calculating steady-state probabilities or finding the probability of being in a particular state at a given time.

• Hidden Markov Models (HMMs): The Forward-Backward algorithm is used for exact inference to compute the posterior probabilities of hidden states given observed data.

- ① Exact Inference
- ② Variable Elimination
- ③ Junction Tree
- ④ Message passing alg.



✓ Approximate inference is used when exact inference is computationally infeasible. It involves finding an approximate solution that is close enough to the exact result. Approximate inference methods trade off some amount of precision for computational efficiency.

using sampling

Crash DA

## Examples of Approximate Inference:

• **Sampling Methods:** Techniques like Markov Chain Monte Carlo (MCMC) methods, including the Metropolis-Hastings algorithm and Gibbs sampling, generate samples from the posterior distribution to approximate it.

• **Variational Inference:** This method approximates a complex probability distribution by a simpler distribution and then minimizes the divergence between the two. The Expectation-Maximization (EM) algorithm and Variational Autoencoders (VAEs) are examples of models using variational inference.

• **Loopy Belief Propagation:** This is an approximate inference technique used in graphical models like Bayesian networks and Markov Random Fields. It iteratively updates beliefs or marginal probabilities but does not guarantee convergence to the exact solution.

• **Expectation Propagation:** A method that approximates the posterior distribution by matching moments between the true distribution and an approximating distribution.



\* Baye's Theorem

$$P(B_i|A) = \frac{P(B_i) \cdot P(A|B_i)}{P(A)}$$



$$P(A, B) = P(A) \cdot P(B|A)$$

$$\rightarrow P(A=0, B=1) = P(A=0) \cdot P(B=1|A=0)$$

$$\rightarrow P(A=0, B=0)$$

$$\rightarrow P(A=1, B=0)$$

$$P(A=1, B=1)$$

Theorem of Total Probability

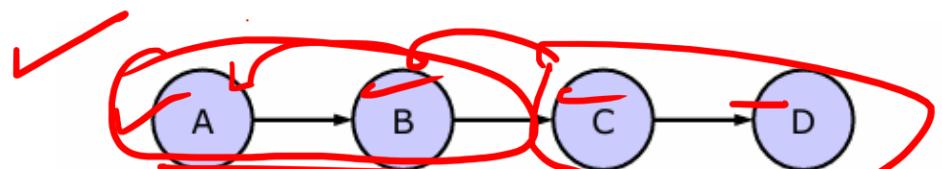
$$P(A) = \sum_{i=1}^n P(B_i) \cdot P(A|B_i)$$

$$P(A) = P(B_0) \cdot P(A|B=0) + P(B=1) \cdot P(A|B=1)$$

$$P(B) = \sum_{i=1}^n P(A_i) \cdot P(B|A_i)$$

$$= P(A=0) \cdot P(B=0|A=0) + P(A=1) \cdot P(B=0|A=1)$$

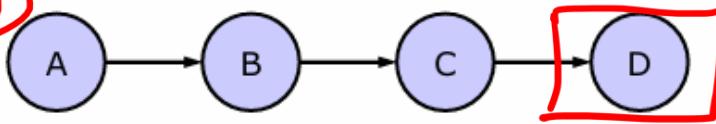
$$+ P(A=0) \cdot P(B=1|A=0) + P(A=1) \cdot P(B=1|A=1)$$



$$\Pr(d) = \sum_{ABC} \Pr(a, b, c, d)$$

$$P(A, B, C, D) = P(A) \cdot P(B|A) \cdot P(C|B) \cdot P(D|C)$$

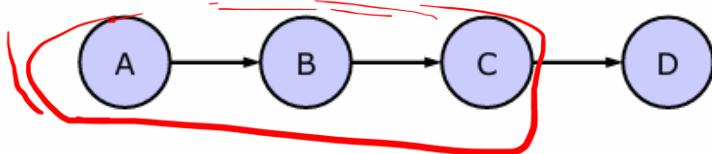
$$\Pr(d) = \Pr(d|c) \Pr(c)$$



$$\Pr(d) = \sum_{ABC} \Pr(a, b, c, d)$$

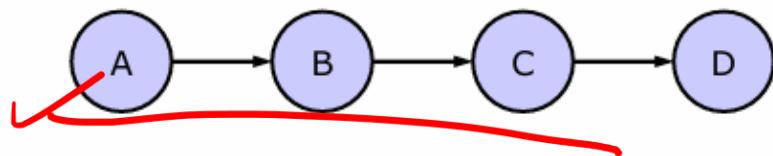
$$= \sum_{ABC} \Pr(d|c) \Pr(c|b) \Pr(b|a) \Pr(a)$$

$$= \sum_C \sum_B \sum_A \Pr(d|c) \Pr(c|b) \Pr(b|a) \Pr(a)$$



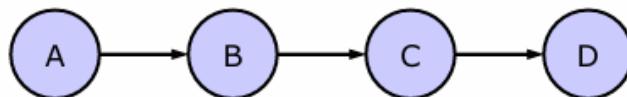
$$\Pr(d) = \sum_{ABC} \Pr(a, b, c, d)$$

$$= \sum_{ABC} \Pr(d|c) \Pr(c|b) \Pr(b|a) \Pr(a)$$



$$\begin{aligned}
 \Pr(d) &= \sum_{abc} \Pr(a,b,c,d) \\
 &= \sum_{abc} \Pr(d|c)\Pr(c|b)\Pr(b|a)\Pr(a) \\
 &= \sum_c \sum_b \sum_a \Pr(d|c)\Pr(c|b)\Pr(b|a)\Pr(a) \\
 &= \sum_c \Pr(d|c) \sum_b \Pr(c|b) \underbrace{\sum_a \Pr(b|a)\Pr(a)}_{\text{Summing out}}
 \end{aligned}$$

# Variable Elimination

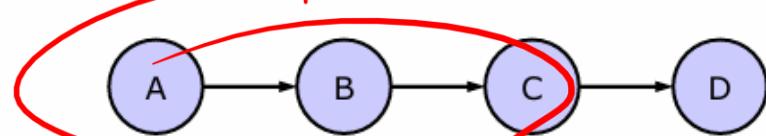


$$\frac{b_1}{b_2} \quad \frac{a_1}{a_2}$$

$$\Pr(d) = \sum_c \Pr(d | c) \sum_B \Pr(c | b) \sum_A \Pr(b | a) \Pr(a)$$

$\Pr(b_1 | a_1) \Pr(a_1) \quad \Pr(b_1 | a_2) \Pr(a_2)$

$\Pr(b_2 | a_1) \Pr(a_1) \quad \Pr(b_2 | a_2) \Pr(a_2)$



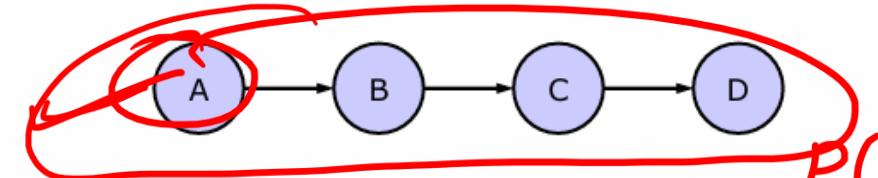
$$\Pr(d) = \sum_c \Pr(d | c) \sum_B \Pr(c | b) \sum_A \Pr(b | a) \Pr(a)$$

$\sum_A \Pr(b_1 | a) \Pr(a)$

$\sum_A \Pr(b_2 | a) \Pr(a)$

$$P(B)$$

$$B=b$$



$$\Pr(d) = \sum_c \Pr(d | c) \sum_B \Pr(c | b) \sum_A \Pr(b | a) \Pr(a)$$

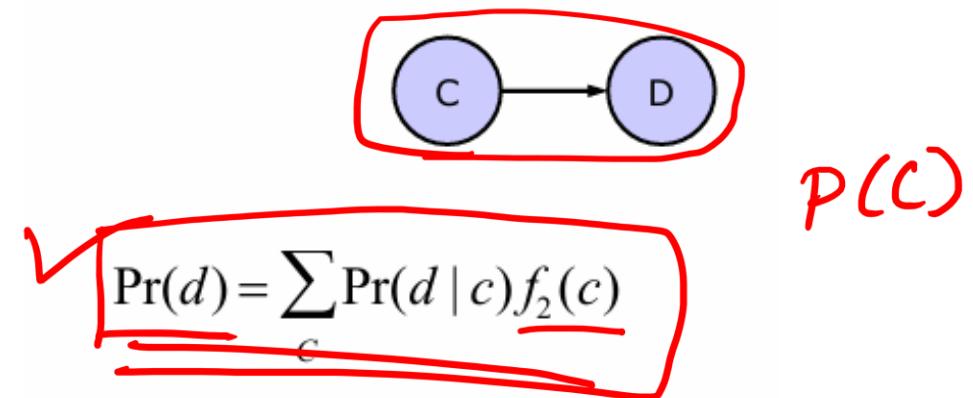
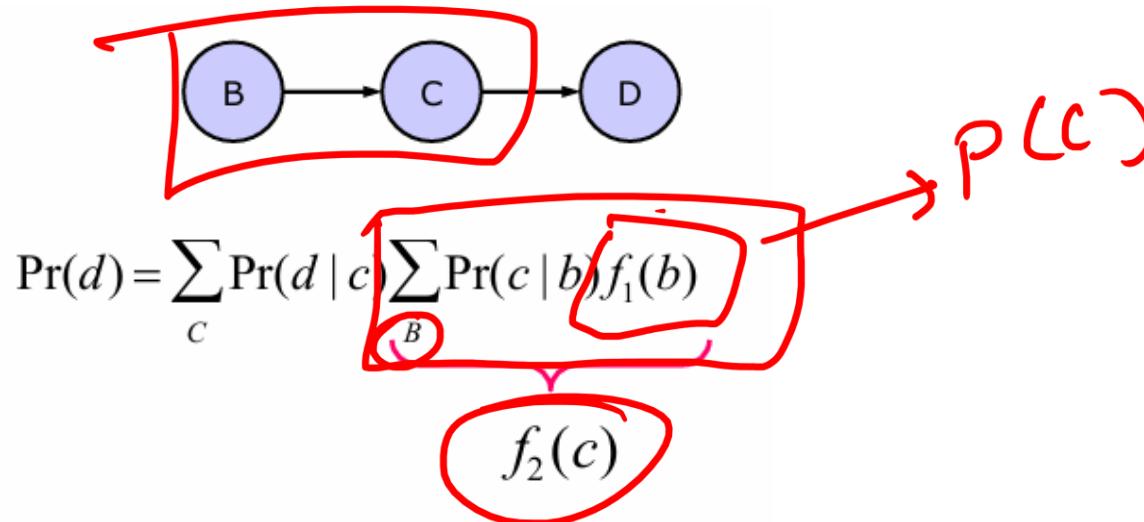
$f_1(b)$

$$f_1(b)$$

$$\underline{f_1(b)}$$

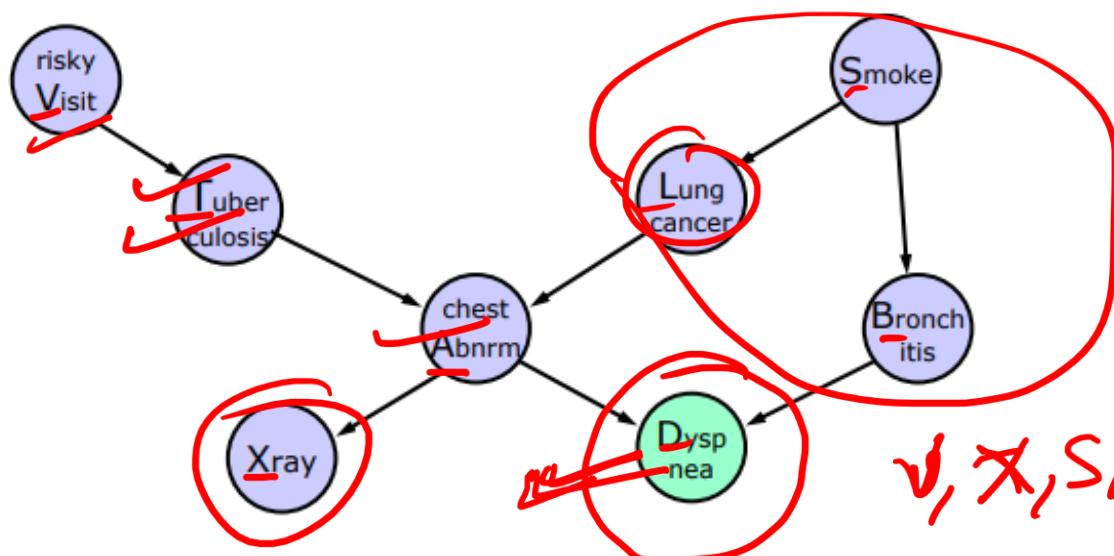
~~$$f_1(b, c)$$~~

~~$$f_1(c) = \sum_B P(c | b) \cdot f_1(b)$$~~

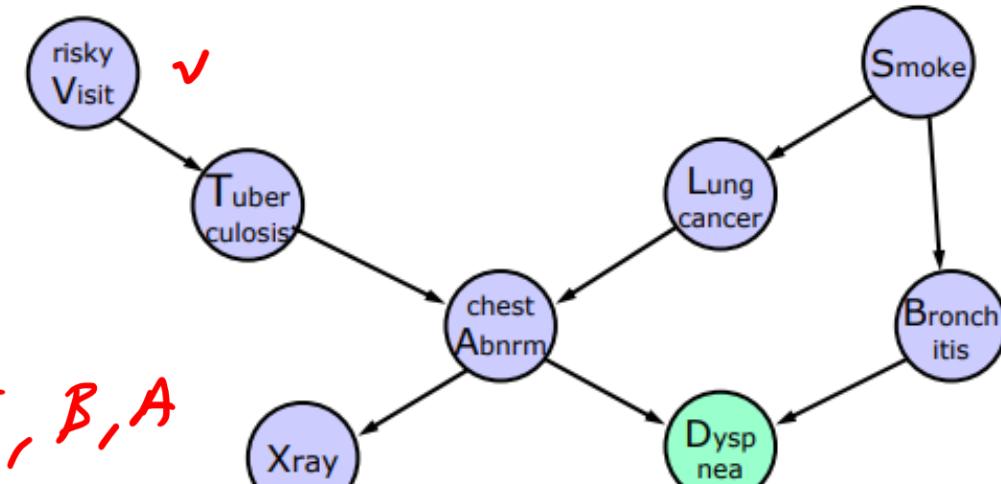


# Variable Elimination Example 2

*Network*



$$P(V, T, X, A, L, S, B, D)$$

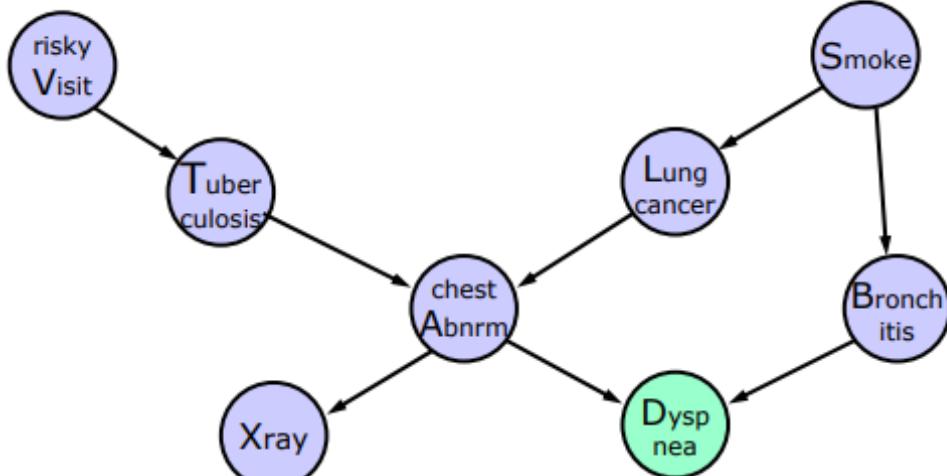


$$\Pr(d) = \sum_{A, B, L, T, S, X, V} \Pr(d | a, b) \Pr(a | t, l) \Pr(b | s) \Pr(l | s) \Pr(s)$$

~~factor  $T \rightarrow f_1(t)$~~

$$P(T)$$

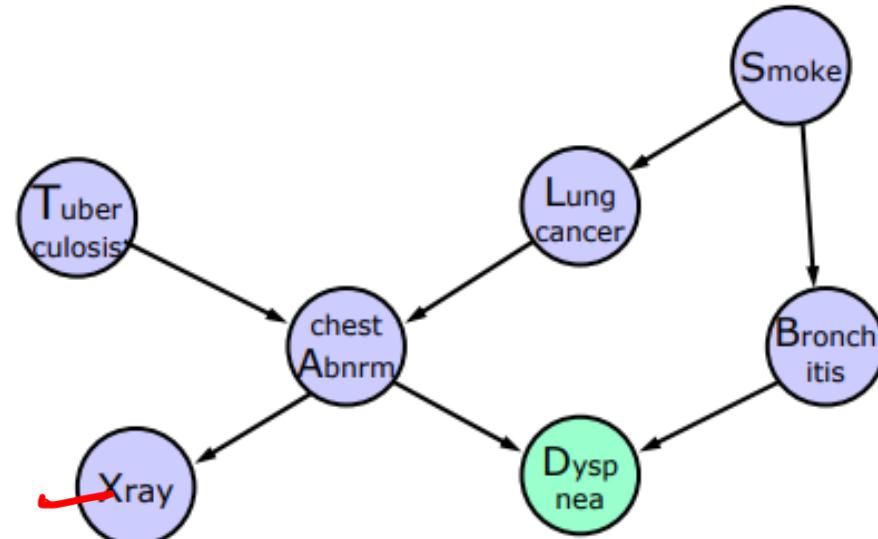
# Variable Elimination Example 2



$$\Pr(d) = \sum_{A,B,L,T,S,X} \Pr(d | a,b) \Pr(a | t,l) \Pr(b | s) \Pr(l | s) \Pr(s)$$

$\Pr(x | a) \sum_v \Pr(t | v) \Pr(v)$   
 $f_1(t)$

$\alpha_1$   
 $\alpha_2$   
 $x_1$   
 $x_2$   
 $CPT$   
 $a_1$   
 $x$



$$\Pr(d) = \sum_{A,B,L,T,S,X} \Pr(d | a,b) \Pr(a | t,l) \Pr(b | s) \Pr(l | s) \Pr(s)$$

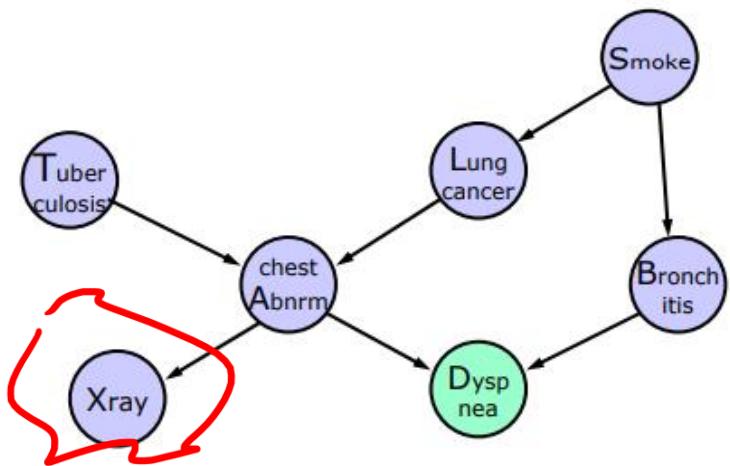
$\Pr(x | a) f_1(t)$

$\sum_x \Pr(x | a) = 1$

# Variable Elimination Example 2

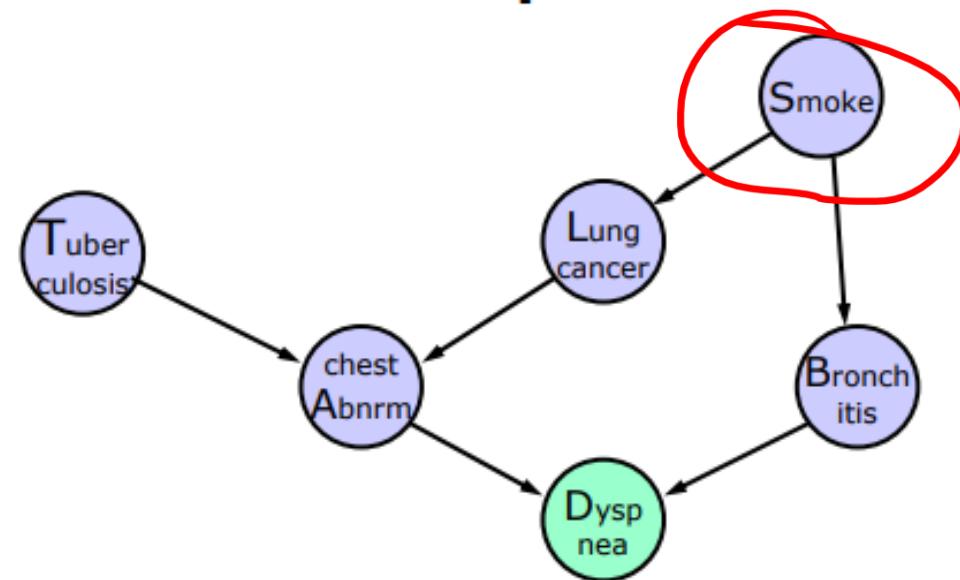


Piyush Wairale



$$\Pr(d) = \sum_{A,B,L,T,S} \Pr(d | a,b) \Pr(a | t,l) \Pr(b | s) \Pr(l | s) \Pr(s) f_1(t)$$

$\sum_{X} \Pr(x | a)$   
1



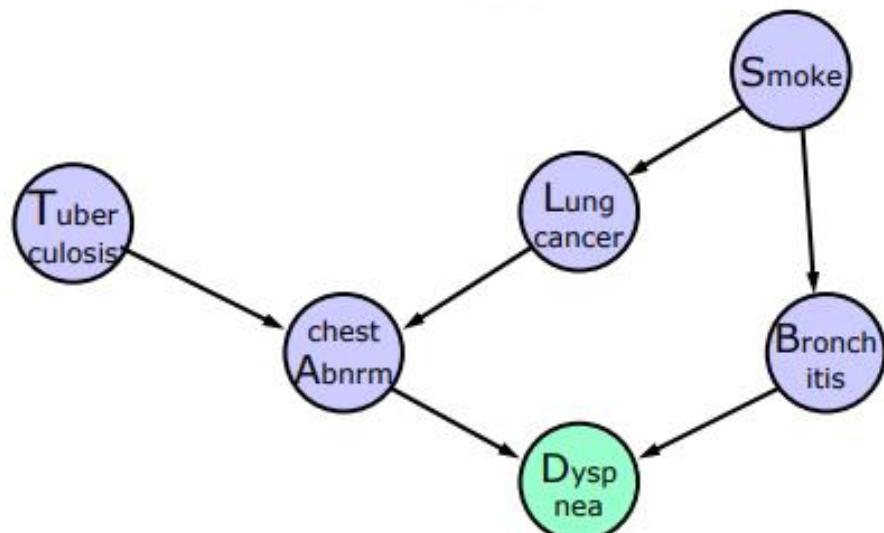
$$\Pr(d) = \sum_{A,B,L,T,S} \Pr(d | a,b) \Pr(a | t,l) \Pr(b | s) \Pr(l | s) \Pr(s) f_1(t)$$

$$\frac{\sum_s \Pr(b | s) \cdot \Pr(l | s) \cdot \Pr(s)}{f_2(b, l)}$$

# Variable Elimination Example 2



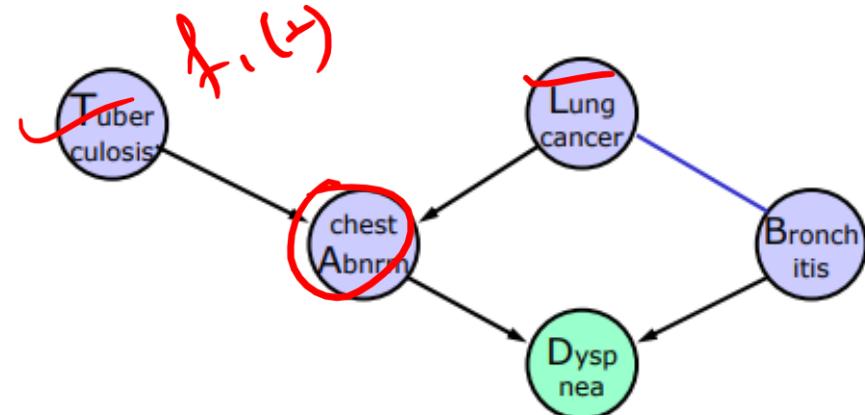
Piyush Wairale



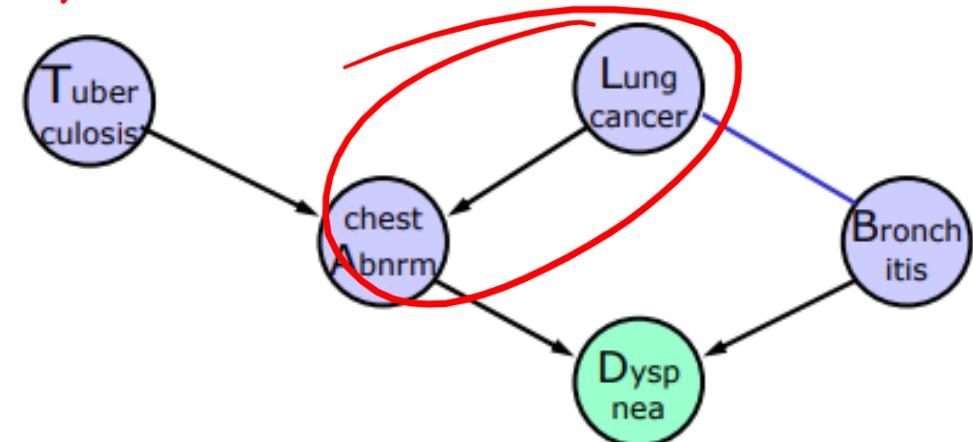
$$\Pr(d) = \sum_{A,B,L,T} \Pr(d | a,b) \Pr(a | t,l) f_1(t) \underbrace{\sum_s \Pr(b | s) \Pr(l | s) \Pr(s)}_{f_1(t)}$$

$$\Pr(d) = \sum_{A,B,L,T} \Pr(d | a,b) \Pr(a | t,l) f_1(t) \underbrace{\sum_s \Pr(b | s) \Pr(l | s) \Pr(s)}_{f_2(b,l)}$$

# Variable Elimination Example 2



$$Pr(a | t, l) \cdot f_1(t)$$



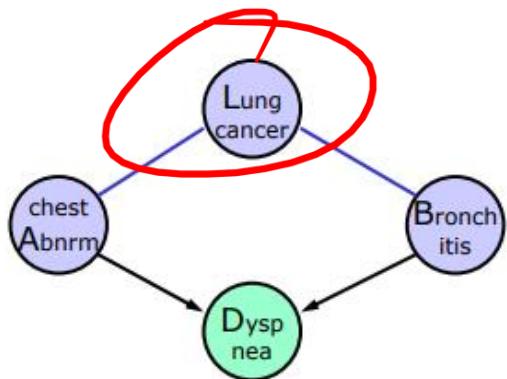
$$Pr(d) = \sum_{A,B,L,T} Pr(d | a,b) \boxed{Pr(a | t,l) f_1(t) f_2(b,l)}$$

$$\checkmark \boxed{f_3(a,l)} \quad \boxed{f_2(b,l)}$$

$$Pr(d) = \sum_{A,B,L} Pr(d | a,b) f_2(b,l) \sum_T \boxed{Pr(a | t,l) f_1(t)}$$

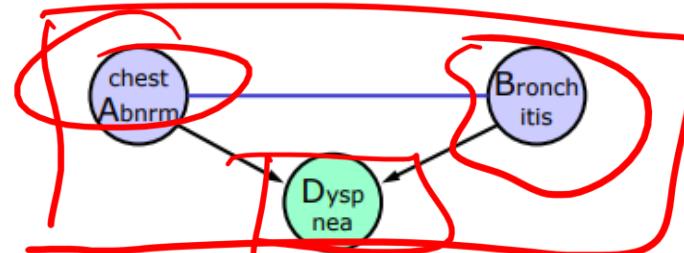
$$\cdot \boxed{f_3(a,l)}$$

# Variable Elimination Example 2



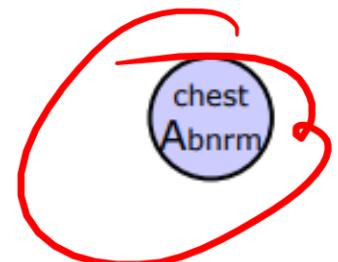
$$\Pr(d) = \sum_{A,B} \Pr(d | a,b) \sum_L f_2(b,l) f_3(a,l)$$

$f_4(a,b)$



$$\Pr(d) = \sum_{A,B} \Pr(d | a,b) f_4(a,b)$$

$f_5(a)$



$$\underline{\Pr(d | a,b) \cdot f_4(a,b)}$$

~~$$\Pr(d) = \sum_a f_5(a)$$~~

# Factors in VE



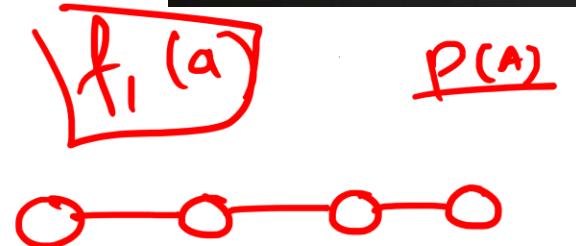
A factor is a function from a set of random variables to a number. Formally, let f denote a factor and let  $X_1$  to  $X_j$  denote the variables in the factor.

A factor is an abstract concept. It can represent a joint probability or a conditional probability. For example, a factor with two variables  $X_1$  and  $X_2$  can represent  $P(X_1 \wedge X_2)$  — this is a joint probability. It can also represent  $P(X_1|X_2)$  — this is a conditional probability. A third possibility is representing  $P(X_1 \text{ and } X_3 = v_3 | X_2)$ . For the third probability, although  $X_3$  appears, it is not a variable since we already assigned a value to it. When you are given a factor, the meanings of the values may not be obvious. For instance, if the factor does not represent a joint distribution, the values may not sum to 1.

For the variable elimination algorithm, we will define a factor for every variable/node in the Bayesian network. The initial factor for each variable/node captures the conditional probability distribution for that variable/node.

$f(E, R)$ :

E	R	val
t	t	0.9
t	f	0.1
f	t	0.0002
f	f	0.9998



$f_1(a, b)$   
 $f(t)$

- ① Restrict
- ② Sum out
- ③ multiply
- ④ Normalise

$$f(E=t, R=t) = 0.9$$
$$f(E=f, R=t) = 0.0002$$

# Factors in VE



Example: Let's look at an example.  $f_1$  is a factor with three variables.

$\text{CPT}$

$f_1(X, Y, Z):$

X	Y	Z	val
t	t	t	0.1
t	t	f	0.9
t	f	t	0.2
t	f	f	0.8
f	t	t	0.4
f	t	f	0.6
f	f	t	0.3
f	f	f	0.7

$X = t$

① Restriction

$$f_1(x=t, Y=f, Z=f) \\ = 0.2$$

$Y = f$

$$f_3(z) = \begin{array}{|c|c|} \hline z & val \\ \hline t & 0.2 \\ f & 0.8 \\ \hline \end{array}$$

$Z = t$

$f_3() = 0.2$

$f_2(Y, Z):$

Y	Z	val
t	t	0.1
t	f	0.9
f	t	0.2
f	f	0.8

# Factors in VE



Example: Let's look at an example of sum out. The factor  $f_1$  has three variables: X, Y, and Z. What happens if we sum out Y? Since we summed out Y, Y will disappear from the factor. The new factor  $f_2$  has X and Z only.

X	Y	Z	val
t	t	t	0.03
t	t	f	0.07
t	f	t	0.54
t	f	f	0.36
f	t	t	0.06
f	t	f	0.14
f	f	t	0.48
f	f	f	0.32

$f_1(X, Y, Z)$ :

② Sum out

addition of val corresponds to  
some  $\times 42$  value  
only

→ for some comb'g  $x \& z$ , we are  
going to add values

X	Z	val
t	t	0.57
t	f	0.43
f	t	0.54
f	f	0.46

$f_2(X, Z)$ :

# Factors in VE



We have two factors.  $f_1$  has variables X and Y.  $f_2$  has variables Y and Z. Let's multiply them together.

$f_1:$

X	Y	val
t	t	0.1
t	f	0.9
f	t	0.2
f	f	0.8

$f_2:$

Y	Z	val
t	t	0.3
t	f	0.7
f	t	0.6
f	f	0.4

$f_3$

③ Multiply

$$f_1 \cup f_2 = XY \cup YZ \\ = \underline{X, Y, Z}$$

Natural Join in SQL

$f_1 \times f_2:$

X	Y	Z	val
t	t	t	0.03
t	t	f	0.07
t	f	t	0.54
t	f	f	0.36
f	t	t	0.06
f	t	f	0.14
f	f	t	0.48
f	f	f	0.32

$0.1 \times 0.3$

$$f_1 \times f_2 \left( \begin{array}{l} X=t, Y=t, \\ Z=f \end{array} \right) \\ = \underline{0.07}$$



**Example:** Suppose that we want to normalize the factor f1. By normalizing, we will produce a new factor f2. Normalize does not change the size of the factor. So f2 has the same size as f1. Each value in f2 is the original value divided by the sum of all the values in f1. The sum is 0.8. So f2 contains 0.25 and 0.75.

④ Normalisation

f1:

Y	val
t	0.2
f	0.6

$$\frac{0.2}{0.2+0.6} = \frac{2}{8}$$

$$= \frac{0.2}{0.8}$$

f2:

Y	val
t	0.25
f	0.75

=

⑤



Approximate inference using sampling refers to a set of methods used to estimate probability distributions, especially in probabilistic models like Bayesian networks, where exact inference is computationally expensive or infeasible.

These methods generate samples from the probability distribution and use these samples to estimate the desired probabilities.

Since the process only approximates the true distribution, it is called approximate inference.

# Approximate Inference

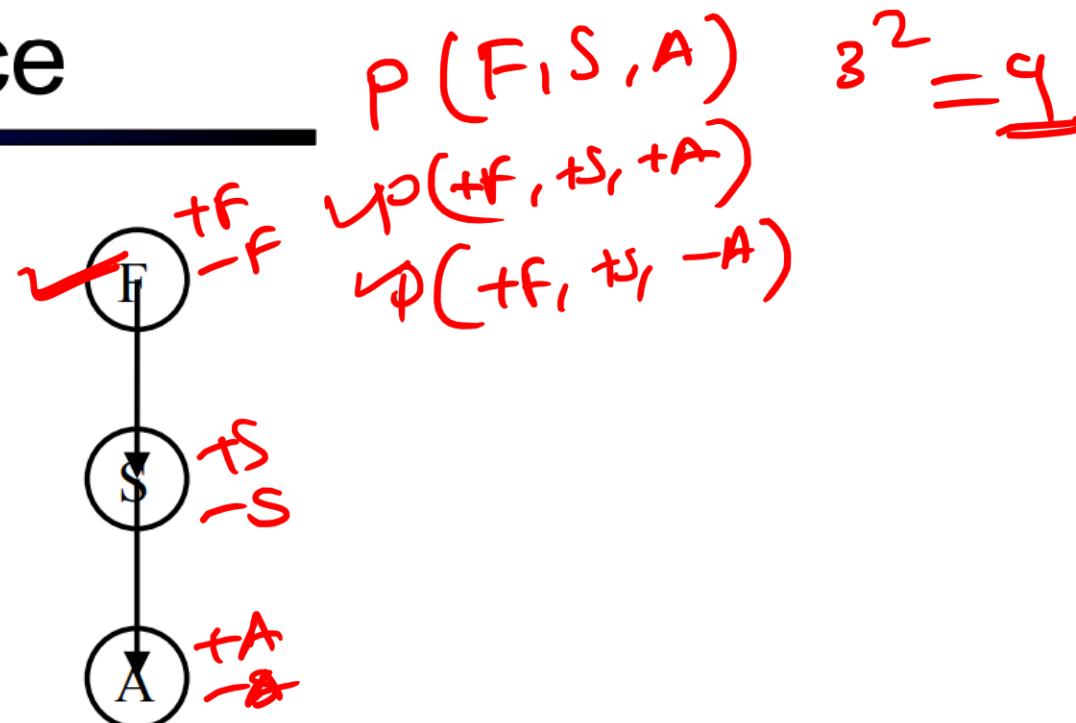
- Simulation has a name: sampling
- Sampling is a hot topic in machine learning, and it's really simple

## Basic idea:

- Draw N samples from a sampling distribution S
- Compute an approximate posterior probability
- Show this converges to the true probability P

## Why sample?

- Learning: get samples from a distribution you don't know
- Inference: getting a sample is faster than computing the right answer (e.g. with variable elimination)





## \* Types of Sampling

- (1) Prior sampling
- (2) Rejection sampling (condition part)
- (3) Likelihood estimate 2024
- (4) Gibbs sampling
- (5) MCMC or Noncov chain monte carlo



Prior sampling is a method used to generate samples from a Bayesian Network without any conditions or evidence. The goal is to sample values for the variables based on the joint probability distribution that the Bayesian network represents.

$P(C | A, B)$

How it works:

1. Start at the root nodes (i.e., nodes with no parents) and sample values for these variables according to their prior probabilities.
2. Move through the network from parents to children, sampling values for each variable based on its conditional probability distribution given the sampled values of its parents.
3. Repeat the process until you have generated the desired number of samples.

CPT

- The generated samples can be used to estimate marginal probabilities by counting occurrences of events.

Prior sampling

# Sampling



Q  
Ans

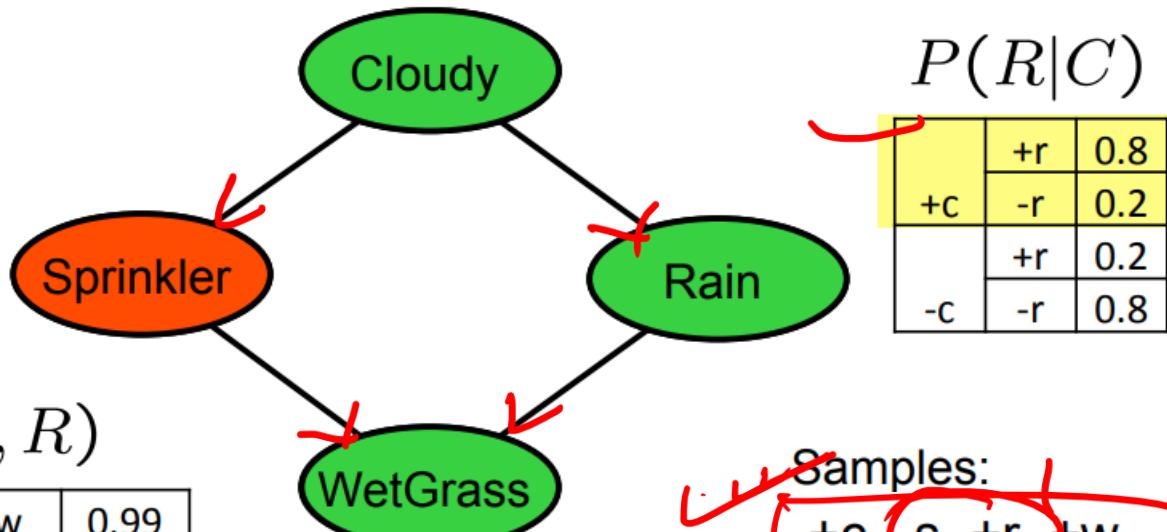
$P(S|C)$

	+s	0.1
+c	-s	0.9
	+s	0.5
-c	-s	0.5

$P(C)$	
+c	0.5
-c	0.5

$$\checkmark P(-s) = ?$$

$$\checkmark P(+w) = ?$$



		Samples:	
		+c, -s, +r, +w	+c, +s, +r, +w
		-c, +s, +r, -w	+c, -s, +r, -w
		+c, -s, +r, +w	+c, -s, +r, -w
		-c, -s, -r, +w	-c, -s, -r, -w

$$P(+c | +s, +r)$$

$$\checkmark P(+c) = P(c = +c)$$

$$= \frac{\text{No. of sample having } +c}{\text{Total No. of sample}}$$

$$= \frac{3}{5}$$

$$= \underline{0.6}$$

$$P(-s, +r) = \frac{2}{5}$$

$$= 0.4$$

$P(W|S, R)$

		+w	0.99
	+r	-w	0.01
		+w	0.90
+s	-r	-w	0.10
		+w	0.90
	+r	-w	0.10
		+w	0.01
-s	-r	-w	0.99



Rejection sampling improves on prior sampling by dealing more effectively with cases where there is evidence (i.e., observed variables). It can be used to estimate conditional probabilities,  $P(X = x | Y = y)$ .

- **Setup:** You have observed some variables (evidence) and want to sample from the conditional distribution given that evidence.
- **Procedure:**
  1. **Generate Prior Samples:** Like in prior sampling, generate samples from the full joint distribution.
  2. **Reject Inconsistent Samples:** After generating a sample, check if it matches the observed evidence (i.e., check if the sample's value for  $Y$  equals the observed value  $y$ ). If the sample is consistent with the evidence, keep it. Otherwise, discard (reject) it.
  3. **Repeat:** Continue generating samples until you have enough that are consistent with the evidence.
  4. **Estimate Probabilities:** Once you have enough samples, estimate probabilities by looking at the proportion of samples that satisfy the desired conditions. For example, to estimate  $P(X = x | Y = y)$ , calculate the fraction of samples where both  $X = x$  and  $Y = y$ , then normalize it to the number of samples where  $Y = y$ .

# Rejection Sampling



5m

$P(S|C)$

	+s	0.1
+c	-s	0.9
	+s	0.5
-c	-s	0.5

$P(C)$

+c	0.5
-c	0.5

$\checkmark +S, +R$

$P(R|C)$

	+r	0.8
+c	-r	0.2
	+r	0.2
-c	-r	0.8

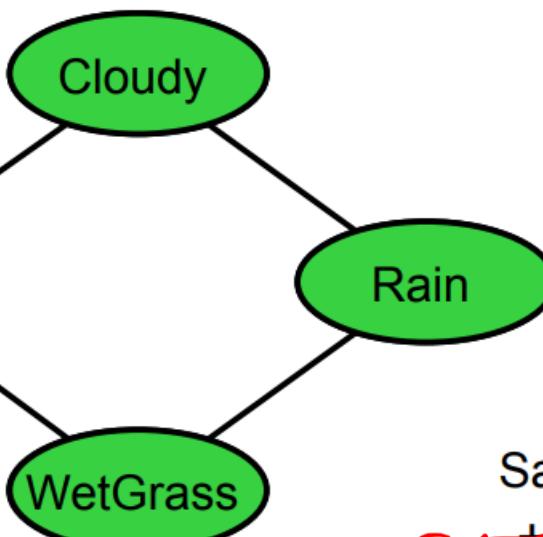
$\checkmark P(+C) / +S + R$

$$= \frac{1}{2}$$

$$= \frac{\text{No. of sample having } +C, +S, +R}{\text{No. of sample having } +S \text{ & } +R}$$

$P(W|S, R)$

		+w	0.99
	+r	-w	0.01
		+w	0.90
+s	-r	-w	0.10
		+w	0.90
	+r	-w	0.10
-s	-r	+w	0.01
		-w	0.99



Samples:

- ~~+c, -s, +r, +w~~
- ~~+c, +s, +r, +w~~
- ~~-c, +s, +r, -w~~
- ~~+c, -s, +r, +w~~
- ~~-c, -s, -r, +w~~

$P(+C|+R, +W)$

$$= ?$$



# Likelihood Weighting

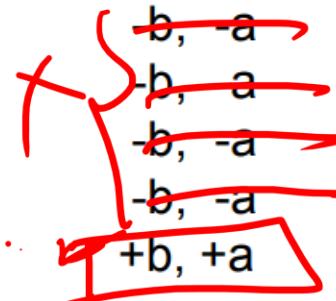
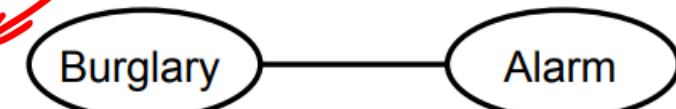
$P(+a | b + c)$  evidence / observed value  
avg

- Problem with rejection sampling:

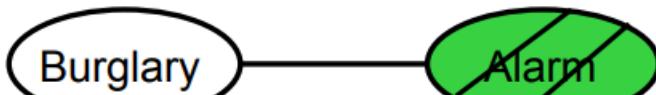
- If evidence is unlikely, you reject a lot of samples
- You don't exploit your evidence as you sample

- Consider  $P(B|+a)$

Query



- Idea: fix evidence variables and sample the rest



+  
 -b +a  
 -b, +a  
 -b, +a  
 -b, +a  
 +b, +a

- Problem: sample distribution not consistent!

- Solution: weight by probability of evidence given parents



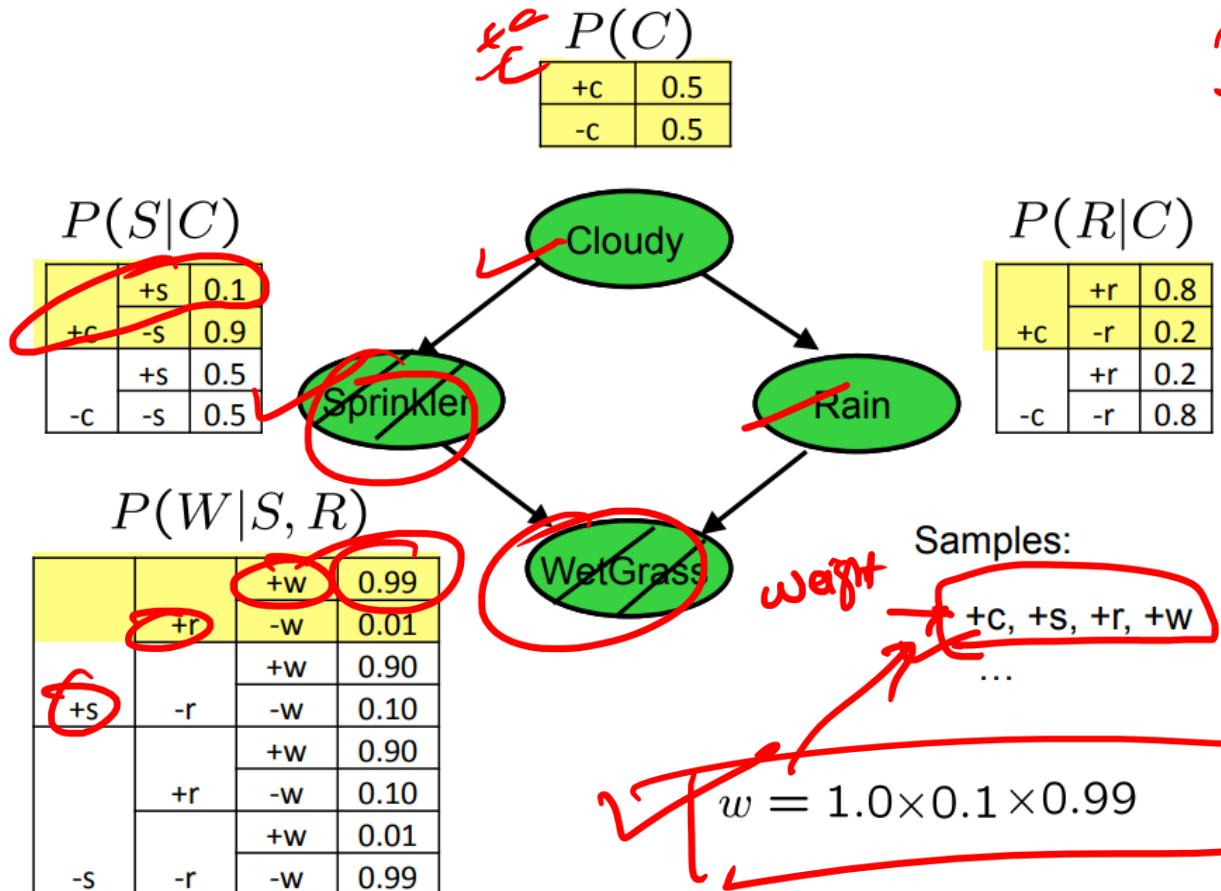
**Likelihood weighting** is a sampling technique used in probabilistic graphical models (such as Bayesian networks) to estimate conditional probabilities.

It is particularly useful when you have evidence (observed variables) and want to avoid the inefficiencies of methods like rejection sampling, which discards many samples.

In rejection sampling, if evidence is rare, many samples will be rejected, making it inefficient. Likelihood weighting addresses this issue by assigning weights to samples based on the likelihood of the evidence, rather than rejecting them.



## Likelihood Weighting



$$\text{evidence} \rightarrow +s, +r$$

$$P(+c | +s, +r)$$

$P(\text{Evidence} | \text{parent})$

$$= \frac{P(+s | +c) \cdot P(+r | +c)}{0.1 \times 0.99} = 0.099$$

samples: +c, +s, +r, +w  
...  
 $w = 1.0 \times 0.1 \times 0.99$

# Likelihood Weighting



$$= P(C) \times P(S|C) \times P(\omega|C)$$

$$\times P(\omega|S, R)$$

$P(C)$	
$+c$	0.5
$-c$	0.5

(JP)

$$\sqrt{P(+c| -s, +\omega)}$$

*evidence*

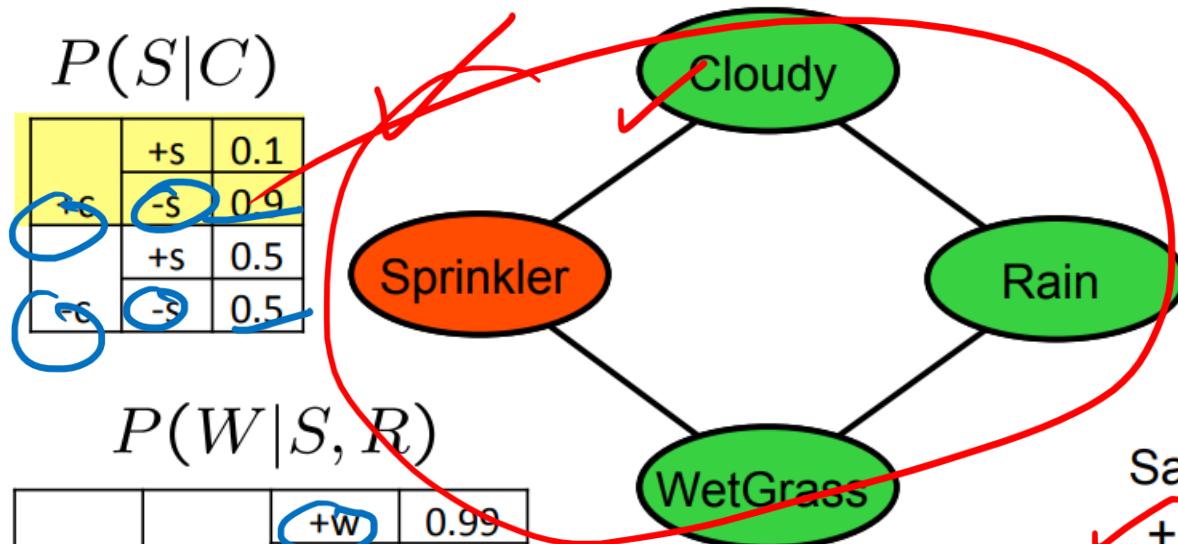
$$\omega_1 = \frac{P(c|c) \cdot P(-s|+c) \cdot P(+\omega|+c)}{P(+\omega|-s, +\omega)}$$

$$\omega_1 = \frac{P(c|c) \cdot P(-s|+c) \cdot P(+\omega|+c)}{P(+\omega|-s, +\omega)}$$

$$= 0.9 \times 0.9 = 0.81$$

$$\omega_2 = \frac{P(-s|-c) \cdot P(+\omega|-s, -c)}{P(+\omega|-s, -c)}$$

$$= 0.5 \times 0.01 = 0.005$$



$P(R|C)$

		$+r$	0.8
$+c$	$-r$	0.2	
$-c$	$+r$	0.2	
$-c$	$-r$	0.8	

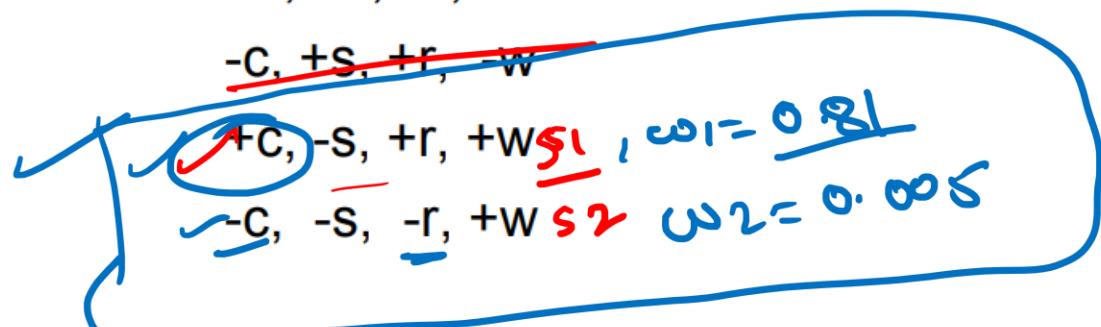
Samples:

$+c, -s, +r, +w$

$+c, +s, +r, +w$

$-c, +s, +r, -w$

$$\sqrt{P(+c| -s, +\omega)} = \frac{0.81}{0.81 + 0.005}$$



# Likelihood Weighting



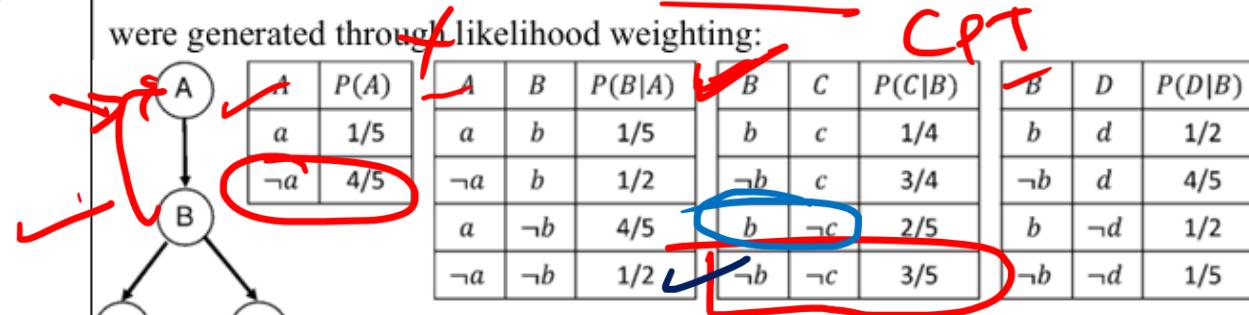
(Q.32)

Consider the Bayes Net containing four Boolean random variables ( $A, B, C, D$ ),

with the following convention:

$$A = \text{True} \Rightarrow A = a, \text{and } A = \text{False} \Rightarrow A = \neg a$$

and similarly for the other variables. The conditional probability tables for the nodes in the network are also indicated in the figure. The following samples were generated through likelihood weighting:



$s_1: (\neg a, \neg b, \neg c, \neg d); s_2: (\neg a, b, \neg c, \neg d); s_3: (\neg a, \neg b, \neg c, d); s_4: (\neg a, b, \neg c, d)$

Estimate the likelihood weight of each sample and thereby estimate

$$P(b|\neg a, \neg c)$$

$$\neg a, \neg c$$

$$P(\text{evident point})$$

$$P(A, B, C, D)$$

$$\equiv P(A) \times P(B|A) \times P(C|B) \times P(D|B)$$

$$w_1 = P(\neg a) \cdot P(\neg b|a)$$

$$s_1: (\neg a, \neg b, \neg c, \neg d)$$

$$w_1 = P(\neg a) \times P(\neg b|\neg a) \times P(\neg c|\neg b) \times P(\neg d|\neg b)$$

$$= \frac{4}{5} \times \frac{3}{5} = \frac{12}{25} = 0.48$$

$$s_2: \neg a, b, \neg c, \neg d.$$

$$w_2 = P(\neg a) \times P(\neg c|b)$$

$$= \frac{4}{5} \times \frac{2}{5} = \frac{8}{25} = 0.32$$

$$s_3: (\neg a, \neg b, \neg c, d)$$

$$w_3 = P(\neg a) \times P(\neg c|\neg b)$$

$$= \frac{4}{5} \times \frac{3}{5} = 0.48$$

$$w_4 = P(\neg a) \cdot P(\neg c|b) = \frac{4}{5} \times \frac{2}{5} = 0.32$$

# Likelihood Weighting



(A)	$s_1: 0.48, s_2: 0.32, s_3: 0.48, s_4: 0.32, P(b \neg a, \neg c) = 0.4$
(B)	$s_1: 0.48, s_2: 0.32, s_3: 0.48, s_4: 0.32, P(b \neg a, \neg c) = 0.64$
(C)	$s_1: 0.32, s_2: 0.48, s_3: 0.48, s_4: 0.32, P(b \neg a, \neg c) = 0.64$
(D)	$s_1: 0.48, s_2: 0.32, s_3: 0.32, s_4: 0.32, P(b \neg a, \neg c) = 0.4$

$$P(b|\neg a, \neg c) = \frac{P(b, \neg a, \neg c)}{\text{sum of wt having } b}$$

query      ev. =      sample  
 summation of wt. of  
 all sample for given  
 evidence

$$= \frac{0.32 + 0.32}{0.48 + 0.32 + 0.48 + 0.32}$$

$$= \frac{0.64}{1.6} = \frac{64}{160}$$

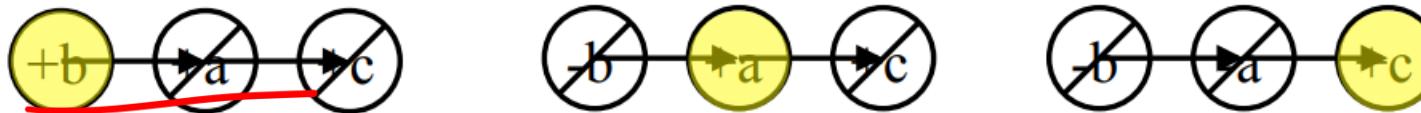
$$= \frac{8}{20} = \frac{4}{10} = \underline{\underline{0.4}}$$

## Markov Chain Monte Carlo\*

- Idea: instead of sampling from scratch, create samples that are each like the last one.

$$p(b|a_t)$$

- Gibbs Sampling: resample one variable at a time, conditioned on the rest, but keep evidence fixed



- Properties: Now samples are not independent (in fact they're nearly identical), but sample averages are still consistent estimators!

- What's the point: both upstream and downstream variables condition on evidence.

$$\begin{array}{c} \xrightarrow{+b, -b} \xrightarrow{+q, -q} \xrightarrow{+c, -c} \\ b \rightarrow a \rightarrow c \end{array}$$

$p(a, b, c) = p(c|b) \cdot p(a|b) \cdot p(b)$

$s_1 = a, -b, c$

$s_2 = a, -b, -c$

$s_3 = a, +b, +c$

$s_1 \vdash +q, -b, c$

$s_2 \vdash +q, +b, c$

# Gibbs Sampling Example

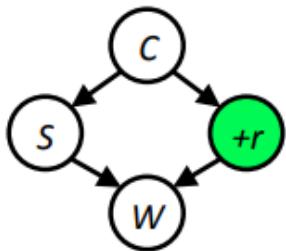
$P(S|+r)$

$+r = \text{evidence}$

*O, Jenahue  
appetizer*

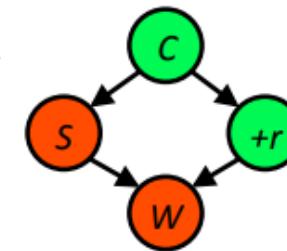
- Step 1: Fix evidence

$R = +r$



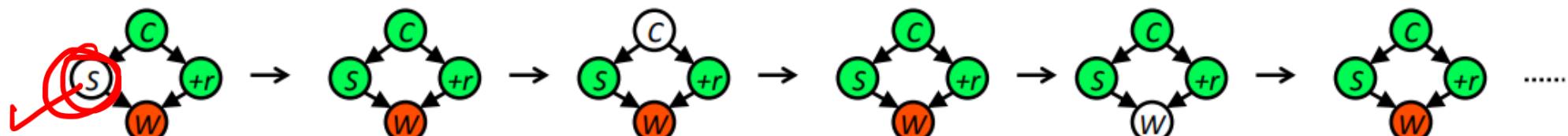
- Step 2: Initialize other variables

Randomly



- Steps 3: Repeat

- Choose a non-evidence variable X
- Resample X from  $P(X | \text{all other variables})$



Sample from  $P(S|+c, -w, +r)$

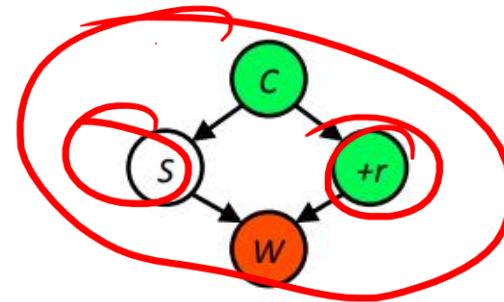
Sample from  $P(C|+s, -w, +r)$

Sample from  $P(W|+s, +c, +r)$

# Sampling One Variable

- Sample from  $P(S|+c, +r, -w)$

$$\begin{aligned}
 P(S|+c, +r, -w) &= \frac{P(S, +c, +r, -w)}{P(+c, +r, -w)} \\
 &= \frac{P(S, +c, +r, -w)}{\sum_s P(s, +c, +r, -w)} \\
 &= \frac{P(+c)P(S|+c)P(+r|+c)P(-w|S,+r)}{\sum_s P(+c)P(s|+c)P(+r|+c)P(-w|s,+r)} \\
 &= \frac{P(+c)P(S|+c)P(+r|+c)P(-w|S,+r)}{P(+c)P(+r|+c)\sum_s P(s|+c)P(-w|s,+r)} \\
 &= \frac{P(S|+c)P(-w|S,+r)}{\sum_s P(s|+c)P(-w|s,+r)}
 \end{aligned}$$



- Many things cancel out – only CPTs with S remain!
- More generally: only CPTs that have resampled variable need to be considered, and joined together



Consider a Bayesian network with three binary variables  $X_1, X_2$ , and  $X_3$ , where each variable depends on the others in some way (e.g., a chain structure). We want to estimate  $P(X_1 = 1 | X_3 = 1)$ .

$$P(X_1 = 1 | X_3 = 1)$$

evolve

Steps:

1. Initialize: Start with arbitrary values for  $X_1, X_2$ , and  $X_3$ , say  $(X_1 = 0, X_2 = 1, X_3 = 1)$ .

2. Iteratively Sample:

- Sample  $X_1$  from its conditional distribution  $P(X_1 | X_2, X_3)$ .
- Sample  $X_2$  from  $P(X_2 | X_1, X_3)$ .
- Sample  $X_3$  from  $P(X_3 | X_1, X_2)$ .

Keep repeating these steps. Since we want to compute the conditional probability given  $X_3 = 1$ , we always fix  $X_3 = 1$  and only update  $X_1$  and  $X_2$ .



100

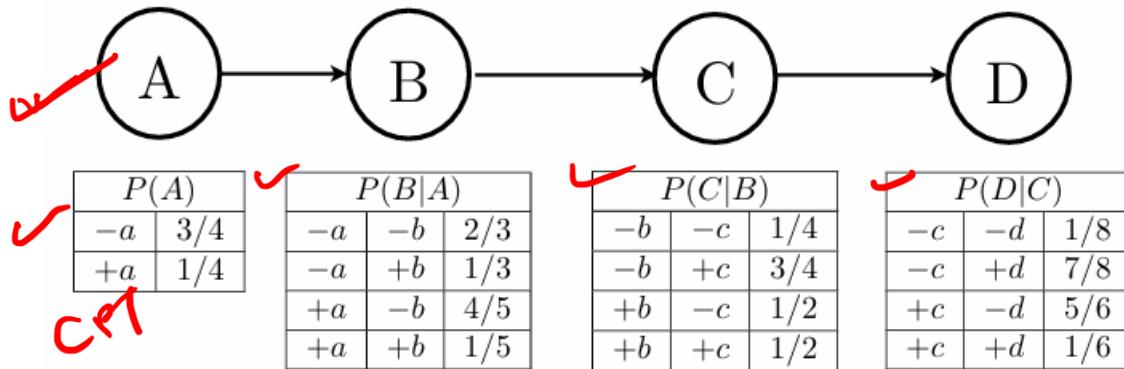
3. Burn-in Period: Discard the first few samples to allow the chain to 'forget' the initial values.
4. Collect Samples: After the burn-in period, collect samples of  $X_1$  and  $X_2$ .
5. Estimate: Use the collected samples to estimate  $P(X_1 = 1 | X_3 = 1)$ .

# Example: Sampling



A

Assume the following Bayes' net, and the corresponding distributions over the variables in the Bayes' net:



jm

$$\underline{P(+c)} = \frac{\text{No. of sample having } +c}{\text{Total no. of samples}} = \frac{5}{8}$$

(a) You are given the following samples:

+a	+b	-c	-d
+a	-b	+c	-d
+a	+b	+c	-d
+a	-b	+c	-d
-b	+c	-d	

(i) Assume that these samples came from performing Prior Sampling, and calculate the sample estimate of  $P(+c)$ .

$5/8$

(ii) Now we will estimate  $P(+c | +a, -d)$ . Above, clearly cross out the samples that would **not** be used when doing Rejection Sampling for this task, and write down the sample estimate of  $P(+c | +a, -d)$  below.

$2/3$

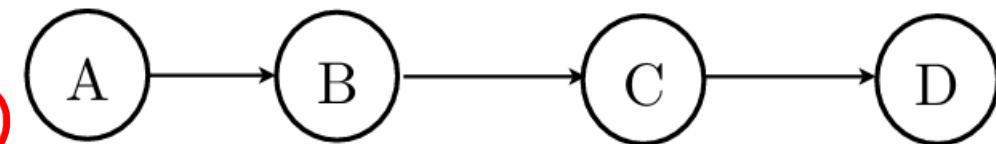
$$\underline{P(+c | +a, -d)} = \frac{\text{No. of sample having } +c}{\text{No. of sample satisfying the given evidence}} = \frac{2}{3}$$

# Example: Sampling



Assume the following Bayes' net, and the corresponding distributions over the variables in the Bayes' net:

$$\begin{aligned} P(D|A) \\ P(A|D) \end{aligned}$$



$P(A)$	
$-a$	$3/4$
$+a$	$1/4$

$P(B A)$		
$-a$	$-b$	$2/3$
$-a$	$+b$	$1/3$
$+a$	$-b$	$4/5$
$+a$	$+b$	$1/5$

$P(C B)$		
$-b$	$-c$	$1/4$
$-b$	$+c$	$3/4$
$+b$	$-c$	$1/2$
$+b$	$+c$	$1/2$

$P(D C)$		
$-c$	$-d$	$1/8$
$-c$	$+d$	$7/8$
$+c$	$-d$	$5/6$
$+c$	$+d$	$1/6$

$$P(a, b, c, d) = P(a) \times P(b|a) \times P(c|b) \times P(d|c)$$

$+b, -d \Rightarrow \text{evidence}$

$P(\text{evidence} | \text{point})$

- (b) Using Likelihood Weighting Sampling to estimate  $P(-a | +b, -d)$ , the following samples were obtained.

Fill in the weight of each sample in the corresponding row.

Sample      Weight

$$-a \quad +b \quad +c \quad -d \quad P(+b | -a)P(-d | +c) = 1/3 * 5/6 = 5/18 = 0.277$$

$$+a \quad +b \quad +c \quad -d \quad P(+b | +a)P(-d | +c) = 1/5 * 5/6 = 5/30 = 1/6 = 0.17$$

$$+a \quad +b \quad -c \quad -d \quad P(+b | +a)P(-d | -c) = 1/5 * 1/8 = 1/40 = 0.025$$

$$-a \quad +b \quad -c \quad -d \quad P(+b | -a)P(-d | -c) = 1/3 * 1/8 = 1/24 = 0.042$$

# Example: Sampling



(c) From the weighted samples in the previous question, estimate  $P(-a | +b, -d)$ .

$$\frac{5/18+1/24}{5/18+5/30+1/40+1/24} = 0.625$$

(d) Which query is better suited for likelihood weighting,  $P(D | A)$  or  $P(A | D)$ ? Justify your answer in one sentence.

$P(D | A)$  is better suited for likelihood weighting sampling, because likelihood weighting conditions only on upstream evidence.

# Example: Sampling



(e) Recall that during Gibbs Sampling, samples are generated through an iterative process.

Assume that the only evidence that is available is  $A = +a$ . Clearly fill in the circle(s) of the sequence(s) below that could have been generated by Gibbs Sampling.

Sequence 1

1 :	+a	-b	-c	+d
2 :	+a	-b	-c	+d
3 :	+a	-b	+c	+d

Sequence 2

1 :	+a	-b	-c	+d
2 :	+a	-b	-c	-d
3 :	-a	-b	-c	+d

Sequence 3

1 :	+a	-b	-c	+d
2 :	+a	-b	-c	-d
3 :	+a	+b	-c	-d

Sequence 4

1 :	+a	-b	-c	+d
2 :	+a	-b	-c	-d
3 :	+a	+b	-c	+a

Gibbs sampling updates one variable at a time and never changes the evidence.

The first and third sequences have at most one variable change per row, and hence could have been generated from Gibbs sampling. In sequence 2, the evidence variables ~~never~~ changed. In sequence 4, the second and third samples have both  $B$  and  $D$  changing.

1  $\rightarrow A = +a \rightarrow$  fixed  
evidence is fixed  
2 one variable at a time