Lecture #32

# Artificial Intelligence for Data Science

**Week-9:**

**Introduction to Probabilistic Reasoning [Part-IV]**

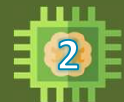Bayesian Network Inference, Parameter Learning, Structure Learning

Course Instructor:

**Dr. Monidipa Das**

**Assistant Professor**

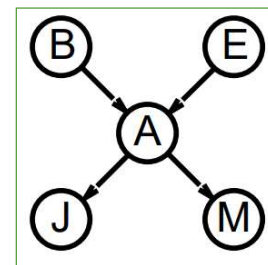**Department of Computational and Data Sciences**

**Indian Institute of Science Education and Research Kolkata, India 741246**

# Inference in Bayesian Networks

- ## Exact inference
  - Inference by enumeration
  - Inference by variable elimination



Simple query :
$$\mathbf{P}(B \mid j, m)$$
$$= \mathbf{P}(B, j, m)/P(j, m)$$
$$= \alpha \mathbf{P}(B, j, m)$$
$$= \alpha \sum_e \sum_a \mathbf{P}(B, e, a, j, m)$$

Rewrite full joint entries using product of CPT entries:
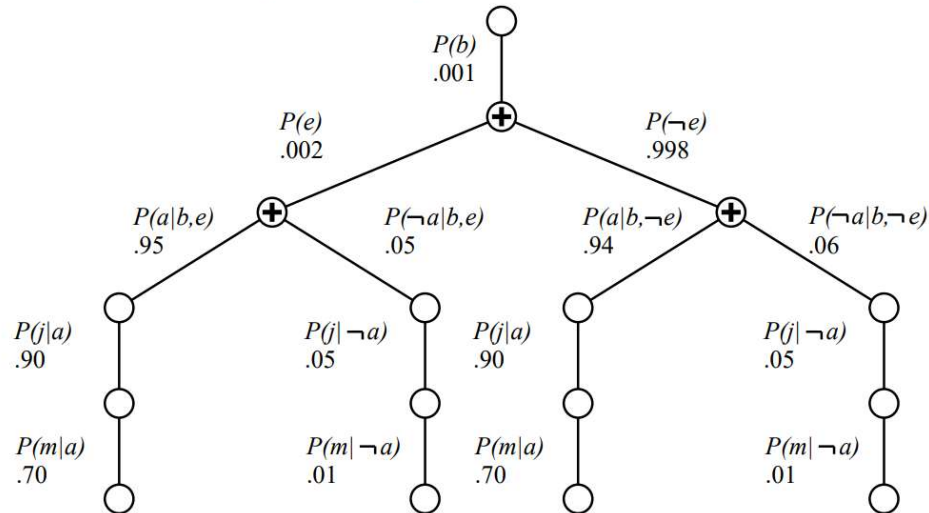$$\mathbf{P}(B \mid j, m)$$
$$= \alpha \sum_e \sum_a \mathbf{P}(B)P(e)\mathbf{P}(a \mid B, e)P(j \mid a)P(m \mid a)$$
$$= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a \mid B, e)P(j \mid a)P(m \mid a)$$

# Evaluation tree

$$\alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a \mid B, e) P(j \mid a) P(m \mid a)$$

$P(b)$
.001

$P(e)$
.002

$P(\neg e)$
.998

$P(a|b,e)$
.95

$P(\neg a|b,e)$
.05

$P(a|b,\neg e)$
.94

$P(\neg a|b,\neg e)$
.06

$P(j|a)$
.90

$P(j|\neg a)$
.05

$P(j|a)$
.90

$P(j|\neg a)$
.05

$P(m|a)$
.70

$P(m|\neg a)$
.01

$P(m|a)$
.70

$P(m|\neg a)$
.01

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Inference in Bayesian Networks

- ## Approximate inference

### Inference by stochastic simulation

Basic idea:

1.  Draw $N$ samples from a sampling distribution $S$

0.5

Coin

2.  Compute an approximate posterior probability $\hat{P}$

3.  Show this converges to the true probability $P$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Inference by stochastic simulation ⑤

**Direct Sampling**

- Basic sampling: sampling with no evidence

- Rejection sampling: reject samples disagreeing with evidence

- Likelihood weighting: use evidence to weight samples

**Markov chain simulation**

- Markov chain Monte Carlo (MCMC): sample from a stochastic process whose stationary distribution is the true posterior

# Approximate inference using MCMC ⑥

- "State" of network = current assignment to all of its variables
- Generate next state by sampling one var. given its Markov Blanket
- Sample each variable in turn, keeping evidence fixed

```
function GIBBS-ASK(X, e, bn, N) returns an estimate of P(X|e)
    local vars: N, a vector of counts for each value of X, initially zero
                Z, the nonevidence variables in bn
                x, the current state of the network, initially copied from e
    initialize x with random values for the variables in Z
    for j = 1 to N do
        for each Z_i in Z do
            set the value of Z_i in x by sampling from P(Z_i | MB(Z_i))
            N[x] ← N[x] + 1 where x is the value of X in x
    return NORMALIZE(N)
```

# MCMC example

To estimate $\mathbf{P}(Rain \mid Sprinkler = true, WetGrass = true)$

1. Apply the Gibbs sampling algorithm with $Sprinkler$ and $WetGrass$ both fixed to $true$
2. Count number of times $Rain$ is $true$ and $false$ in the samples

**Example:**
Visit 100 states; 31 have $Rain = true$, 69 have $Rain = false$

$\hat{\mathbf{P}}(Rain \mid Sprinkler = true, WetGrass = true)$
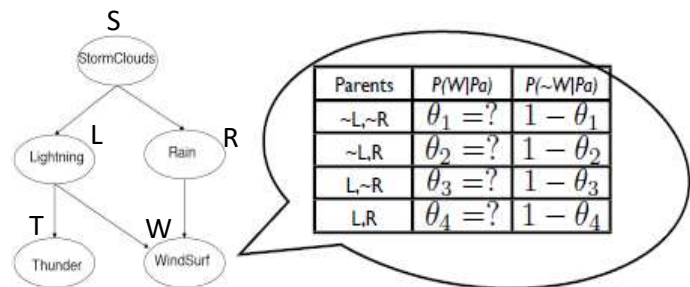$= \textsc{Normalize}(\langle 31, 69 \rangle) = \langle 0.31, 0.69 \rangle$

$P(C) = .5$

Cloudy

| C | P(S) |
|---|------|
| T | .10 |
| F | .50 |

Sprinkler

Rain

| C | P(R) |
|---|------|
| T | .80 |
| F | .20 |

Wet Grass

| S | R | P(W) |
|---|---|------|
| T | T | .99 |
| T | F | .90 |
| F | T | .90 |
| F | F | .00 |

# Learning Parameters

- Parameters ($\theta$)
  - Probabilities in the CPTs for all the variables in the network

- Learning parameters
  - Infer $\theta$ from data, given G

S

StormClouds

L

Lightning

R

Rain

T

Thunder

W

WindSurf

| Parents | P(W\|Pa) | P(~W\|Pa) |
|---------|----------|-----------|
| ~L,~R | $\theta_1 = ?$ | $1 - \theta_1$ |
| ~L,R | $\theta_2 = ?$ | $1 - \theta_2$ |
| L,~R | $\theta_3 = ?$ | $1 - \theta_3$ |
| L,R | $\theta_4 = ?$ | $1 - \theta_4$ |

# Learning Parameters

- *G* is a given DAG over *N* variables

- Goal: Estimate $\theta$ from i.i.d data $D = X = (x_1, \dots, x_M)$,
  where $M$ is the number of records
  Each record: $x_m = \{x_m^1, x_m^2, \dots, x_m^N\}$

- Complete Observability (no missing values)

# Bayes' formula as touchstone

$$prob(\Theta|\mathcal{X}) \quad = \quad \frac{prob(\mathcal{X}|\Theta) \cdot prob(\Theta)}{prob(\mathcal{X})}$$

$$posterior \quad = \quad \frac{likelihood \cdot prior}{evidence}$$

# Parameter Estimation Methods

11

```
Parameter Estimation
Methods
    ├── Frequentist Method
    │       └── E.g.: MLE
    └── Bayesian Method
            └── E.g.: MAP, Bayesian Estimation
```

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Frequentist vs. Bayesian

12

| Bayesian Method | Frequentist Method |
|---|---|
| Uses probabilities for both hypotheses and data | Never uses or gives the probability of a hypothesis (no prior or posterior) |
| Depends on the prior and likelihood of observed data | Depends on the likelihood for both observed and unobserved data. |
| Requires one to know or construct a 'subjective prior' | Does not require a prior. |
| Dominated statistical practice before the 20th century | Dominated statistical practice during the 20th century. |
| May be computationally intensive due to integration over many parameters | Tends to be less computationally intensive. |

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# MLE: Maximum Likelihood Estimator  `13`

- We seek that value for Θ which maximizes the likelihood (i.e. $prob(X|\Theta)$)
- We denote such a value of Θ by $\widehat{\Theta}_{ML}$
- $\widehat{\Theta}_{ML}$ maximizes

$$\prod_{\mathbf{x}_i \in \mathcal{X}} prob(\mathbf{x}_i|\Theta)$$

  (assuming independent observations)

- Simply, $\widehat{\Theta}_{ML}$ maximizes

$$\mathcal{L} = \sum_{\mathbf{x}_i \in \mathcal{X}} log\ prob(\mathbf{x}_i|\Theta)$$

- Hence,

$$\widehat{\Theta}_{ML} = \underset{\Theta}{argmax}\ \mathcal{L} \qquad \frac{\partial \mathcal{L}}{\partial \theta_i} = 0 \qquad \forall\ \theta_i \in \Theta$$

# MAP: Maximum A Posteriori estimate  `14`

- We seek that value for Θ which maximizes the posterior (i.e. $prob(\Theta|X)$)
- We denote such a value of Θ by $\widehat{\Theta}_{MAP}$

$$\widehat{\Theta}_{MAP} = \underset{\Theta}{argmax}\ prob(\Theta|\mathcal{X}) \qquad \boxed{\widehat{\Theta}_{MAP} = \underset{\Theta}{argmax}\left(\sum_{\mathbf{x}_i \in \mathcal{X}} log\ prob(\mathbf{x}_i|\Theta) + log\ prob(\Theta)\right)}$$

$$= \underset{\Theta}{argmax}\ \frac{prob(\mathcal{X}|\Theta)\ \cdot\ prob(\Theta)}{prob(\mathcal{X})}$$

$$= \underset{\Theta}{argmax}\ prob(\mathcal{X}|\Theta)\ \cdot\ prob(\Theta)$$

$$= \underset{\Theta}{argmax}\ \prod_{\mathbf{x}_i \in \mathcal{X}} prob(\mathbf{x}_i|\Theta)\ \cdot\ prob(\Theta)$$

# Bayesian Estimation

15

- $\widehat{\Theta}_B$ maximizes full posterior

$$prob(\Theta|\mathcal{X}) \quad = \quad \frac{prob(\mathcal{X}|\Theta) \;\cdot\; prob(\Theta)}{prob(\mathcal{X})}$$

- The denominator in the Bayes' Rule cannot be ignored.

- The denominator, known as the probability of evidence, is related to the other probabilities that make their appearance in the Bayes' Rule by

$$prob(\mathcal{X}) \;=\; \int_{\Theta} prob(\mathcal{X}|\Theta) \cdot prob(\Theta) \; d\Theta$$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# MLE vs. MAP vs. Bayesian Estimation

16

- MLE does NOT allow us to inject our prior beliefs about the likely values for $\Theta$ in the estimation calculations.
- MAP allows for the fact that the parameter vector $\Theta$ can take values from a distribution that expresses our prior beliefs regarding the parameters.
- Both MLE and MAP return only single and specific values for the parameter $\Theta$.
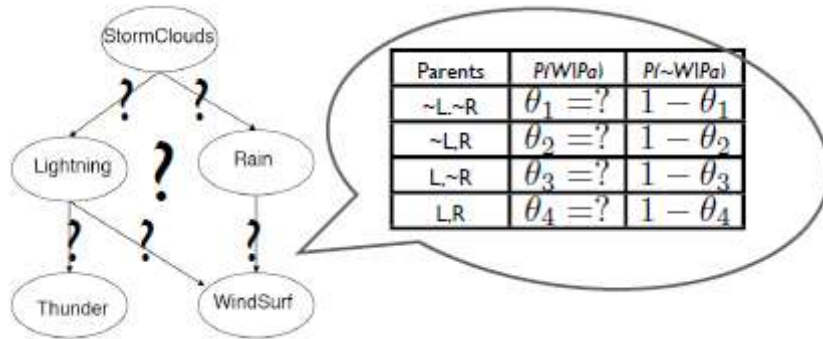- Bayesian estimation, by contrast, calculates fully the posterior distribution $prob(\Theta|X)$.

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Learning Graph Structure

- Structure Learning:
  - inferring $G$ and $\theta$ from data
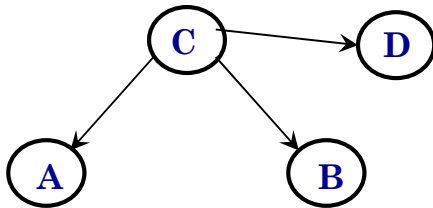
# Structural Learning Methods

- Constraint Based
  - Test independencies
  - Add edges according to the tests

- Search and Score
  - Define a selection criterion that measures goodness of a model
  - Search in the space of all models (or orders)

- Mix models (recent)
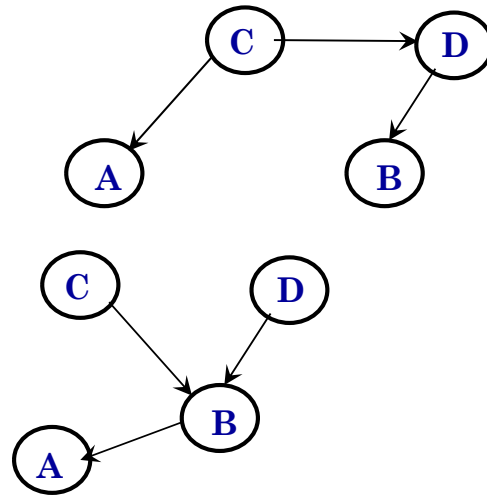  - Test for almost all independencies
  - Search and scoring

# Search and Score

- Learning from the data itself
  - Space of possible structures
  - Scoring each structure

$$P(\Theta|X) \propto P(\Theta).P(X|\Theta)$$

# Search and Score

Find $\hat{\mathbf{G}} = \arg\max_{G} \mathbf{Score(G)}$

- **Heuristic search:**
  - Greedy local search
  - Best-first search
  - Simulated annealing

Add C->B

Delete S->B

Reverse S->B

# Bayesian Score

- Main principle of the Bayesian approach

Marginal likelihood     Prior over structures

$$P(G \mid D) = \frac{P(D \mid G)P(G)}{P(D)}$$

Marginal probability of Data

P(D) does not depend on the network

**Bayesian Score:** $Score_B(G:D) = \log P(D \mid G) + \log P(G)$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

---

# Questions?

Dr. Monidipa Das, Department of CDS, IISER Kolkata