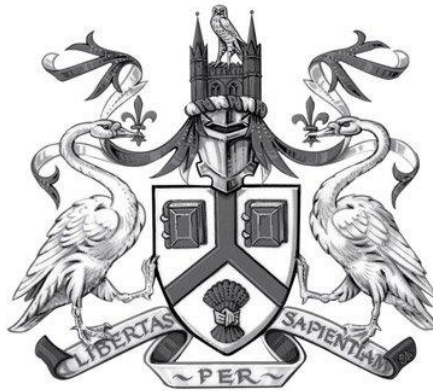


Access Code: **489141**

**CMP9794M**

# **Advanced Artificial Intelligence**

[Heriberto Cuayahuitl](#)



UNIVERSITY OF  
**LINCOLN**

School of Engineering and Physical Sciences

Thinking humanly	Thinking rationally
Acting humanly	Acting rationally

# Last Week

Fully-observable vs. partially observable

Single-agent vs. multi-agent

Deterministic vs. stochastic

Episodic vs. sequential

Static vs. dynamic

Discrete vs. continuous

Known vs. unknown

- Main approaches to AI
- Agents & environments
- History and developments
- Probability theory
- Naïve Bayes classifier

$$P(A \mid B) = P(A \wedge B) / P(B)$$

$$P(A \wedge B) = P(A \mid B) * P(B)$$

$$P(A \mid B) + P(\neg A \mid B) = 1$$

$$P(B) = \sum_a P(A=a, B) = \sum_a P(A=a \mid B) * P(B)$$

$$P(A) + P(\neg A) = 1, \text{ therefore } P(\neg A) = 1 - P(A)$$

$$P(B) + P(\neg B) = 1, \text{ therefore } P(B) = 1 - P(\neg B)$$

$$P(A \wedge B) = P(B \wedge A)$$

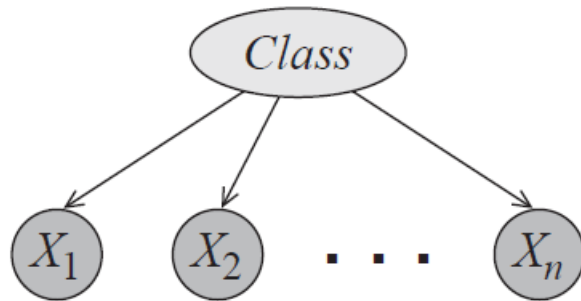
$$P(A \mid B) \neq P(B \mid A)$$

$$P(A \mid B) = (P(B \mid A) * P(A)) / P(B)$$

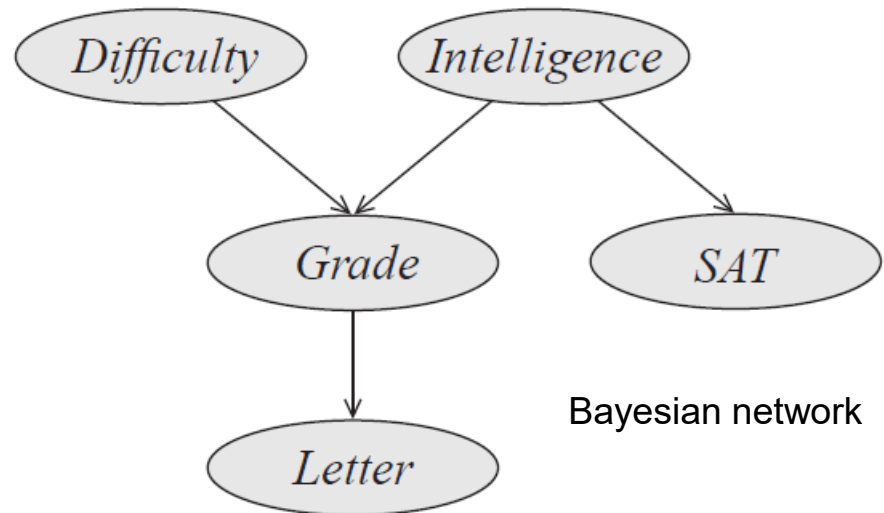
$$Y = \operatorname{argmax}_{y_k} P(Y = y_k) \prod_i P(X_i \mid Y = y_k)$$

# From Naïve Bayes to Bayesian Nets

Naïve bayes is a simple Bayesian Network (BN) with a strong independence assumption, which is relaxed in BNs via not so simple structures.

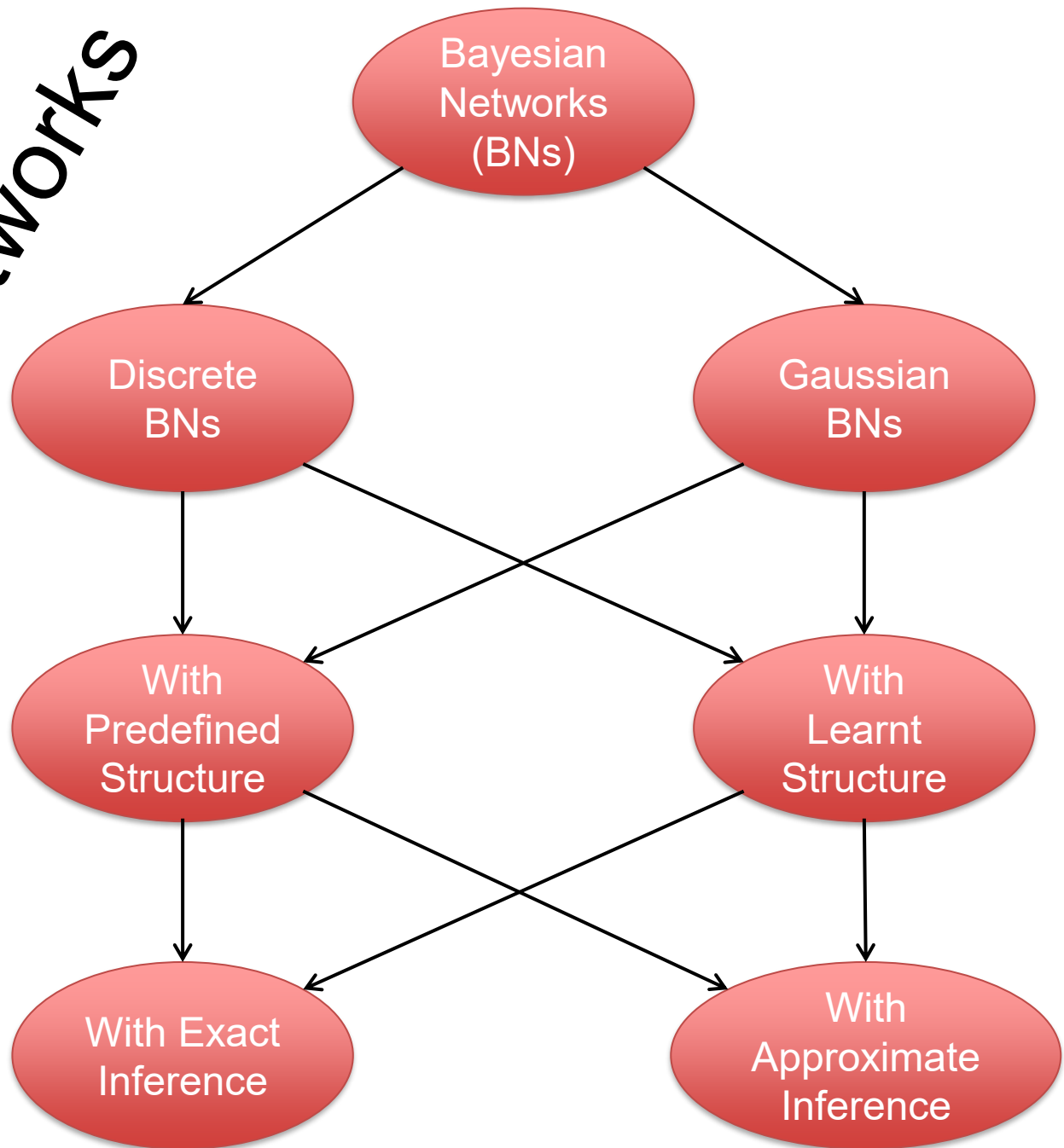


Naïve Bayes graphical model



Bayesian network

# Taxonomy of Bayesian Networks



# Today

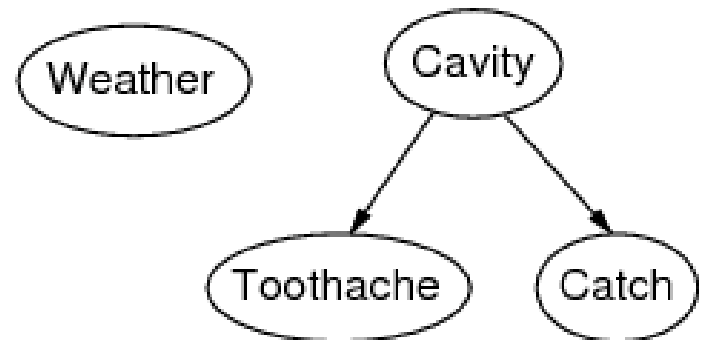
- **Introduction to Discrete Bayesian networks**
  - Graphical and probabilistic representation
  - Parameter learning
- Algorithms for exact inference
  - Inference by enumeration
  - Inference by variable elimination

# Bayesian Networks

- **Bayesian Networks** (Bayes Nets or Belief Nets) can represent any full joint probability distribution—and they can do so very concisely!
- **Syntax:**
  - a set of nodes, one per random variable
  - a directed acyclic graph (link=“directly influences”)
  - a conditional distribution for each node given its parents:  $P(X_i | \text{parents}(X_i))$

# Bayesian Networks (BNs)

- Each node of a BN is represented by a **conditional probability table (CPT)**—a probability distribution over  $X_i$  for each combination of parent values.
- The topology of a network encodes conditional independence assertions:
  - *Weather* is independent of the other variables
  - *Toothache* and *Catch* are conditionally independent given *Cavity*



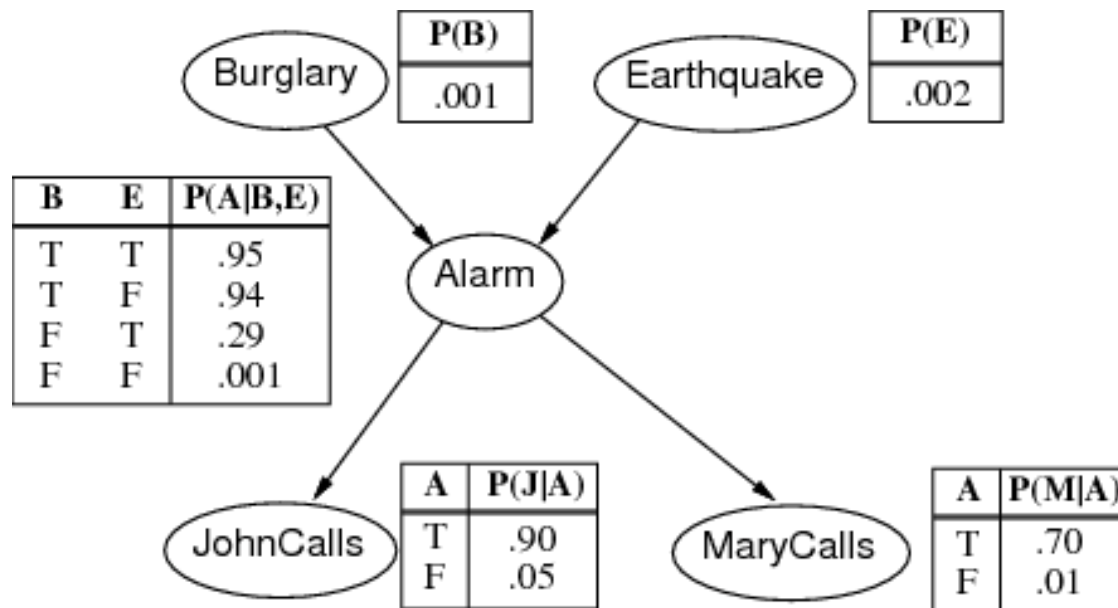
# Example Scenario

- Excerpt from Russell and Norvig (2016) *“I am at work, my neighbour John calls to say my alarm is ringing, and my neighbour Mary doesn't call. Sometimes the alarm is set off by minor earthquakes. Is there a burglar?”*
- Random variables (binary):
  - B=Burglar
  - E=Earthquake
  - A=Alarm
  - J=JohnCalls
  - M=MaryCalls



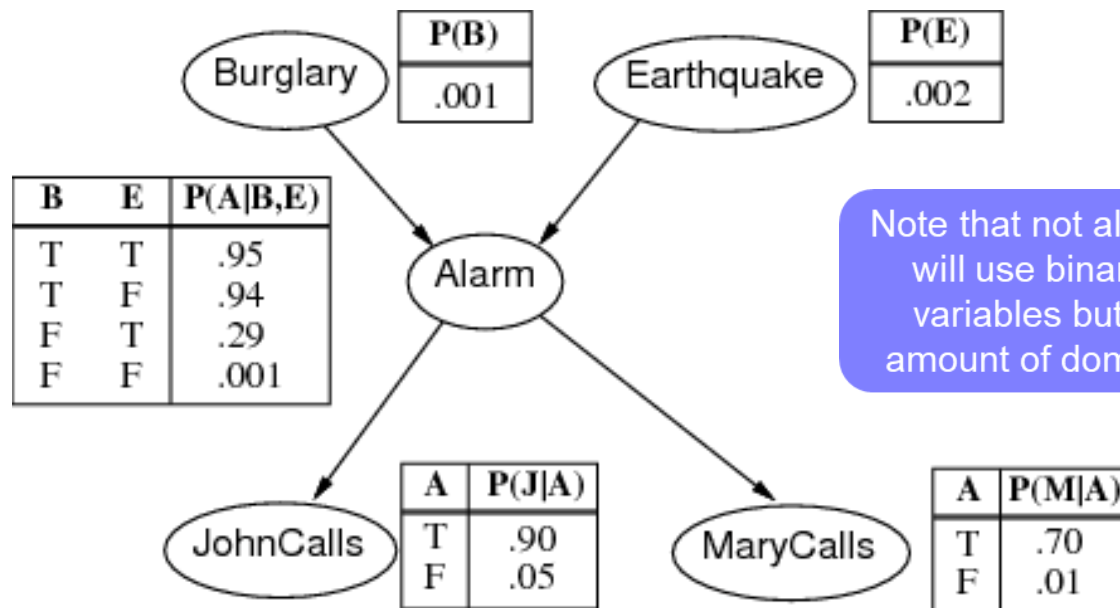
# Example Scenario

- The network topology reflects “causal” knowledge:
  - A burglar can set the alarm on
  - An earthquake can set the alarm on
  - The alarm can cause Mary to call
  - The alarm can cause John to call



# Example Scenario

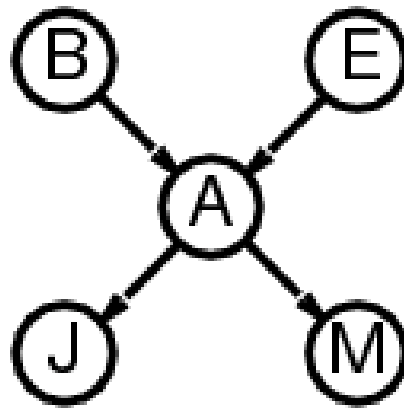
- The network topology reflects “causal” knowledge:
  - A burglar can set the alarm on
  - An earthquake can set the alarm on
  - The alarm can cause Mary to call
  - The alarm can cause John to call



Note that not all Bayes nets will use binary random variables but a varying amount of domain values.

# Compactness

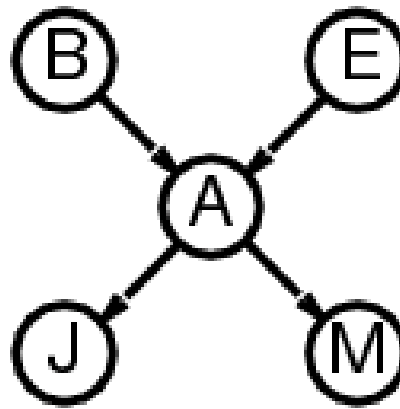
- A CPT for binary variable  $X_i$  with  $k$  binary parents has  $2^k$  rows for the combinations of parent values



- Each row requires one real number  $p$  for  $X_i = \text{true}$ , and one for  $X_i = \text{false}$  (i.e.,  $\neg p = 1 - p$ ). For the burglary net,  $2^0 + 2^0 + 2^2 + 2^1 + 2^1 = 10$  numbers.

# Compactness

- A CPT for binary variable  $X_i$  with  $k$  binary parents has  $2^k$  rows for the combinations of parent values



- Each row requires one real number  $p$  for  $X_i = \text{true}$ , and one for  $X_i = \text{false}$  (i.e.,  $\neg p = 1 - p$ ). For the burglary net,  $2^0 + 2^0 + 2^2 + 2^1 + 2^1 = 10$  numbers. Full enumeration requires  $2^1 + 2^1 + 2^3 + 2^2 + 2^2 = 20$

# Compactness

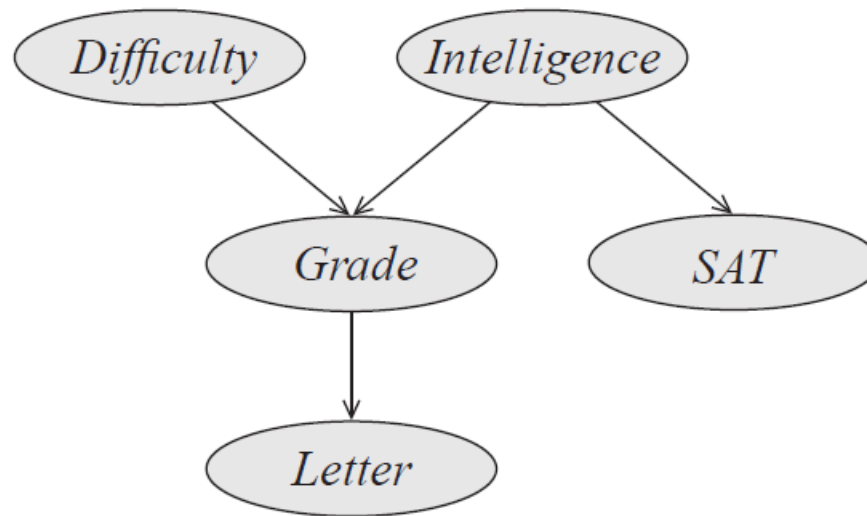
- If each variable has no more than  $k$  parents, the complete network requires  $n * 2^k$  numbers.
- **[Question]** What is the number of probabilities in a Bayesian Network with 30 binary random variables, each with 5 parents – using compact enumeration?
- **[Question]** What is the number of probabilities in the full joint distribution – using full enumeration)?

# Compactness

- If each variable has no more than  $k$  parents, the complete network requires  $n * 2^k$  numbers.
- **[Question]** What is the number of probabilities in a Bayesian Network with 30 binary random variables, each with 5 parents – using compact enumeration?  
 $n * 2^k = 30 * 2^5 = 960$
- **[Question]** What is the number of probabilities in the full joint distribution – using full enumeration)?  
 $2^n = 2^{30} = 1,073,741,824$

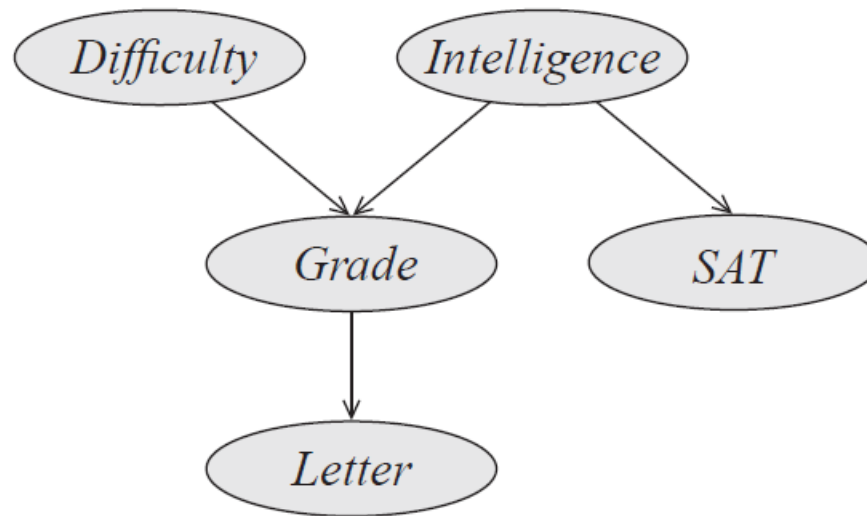
# Number of Probabilities in Bayes Nets

**[Question]** How many probabilities are required by all CPTs of the Bayesian Network below considering that all variables except  $G$  are binary— $G$ 's domain size is 3?



# Number of Probabilities in Bayes Nets

**[Question]** How many probabilities are required by all CPTs of the Bayesian Network below considering that all variables except  $G$  are binary— $G$ 's domain size is 3?

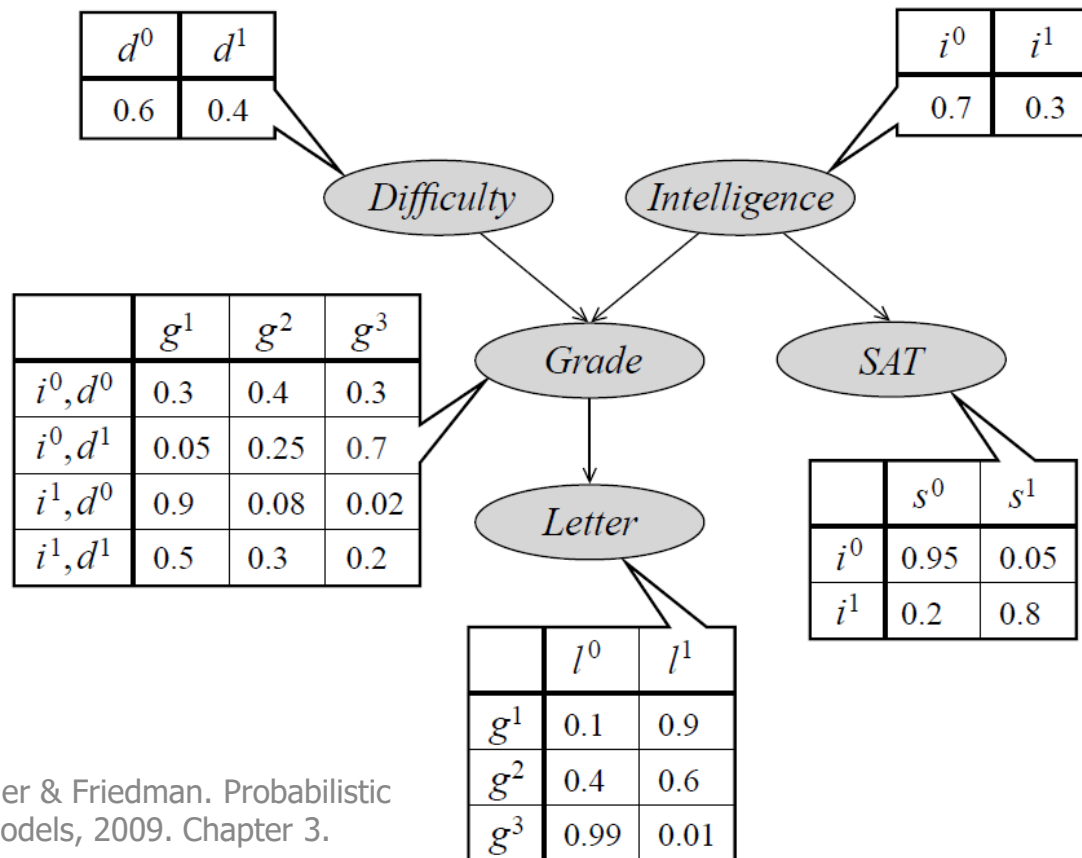


The answer is  $2+2+12+4+6=26$  due to  $|D| = 2$ ,  $|I| = 2$ ,  $|G| = 3 * 2 * 2 = 12$ ,  $|SAT| = 2 * 2 = 4$ ,  $|L| = 3 * 2 = 6$ .



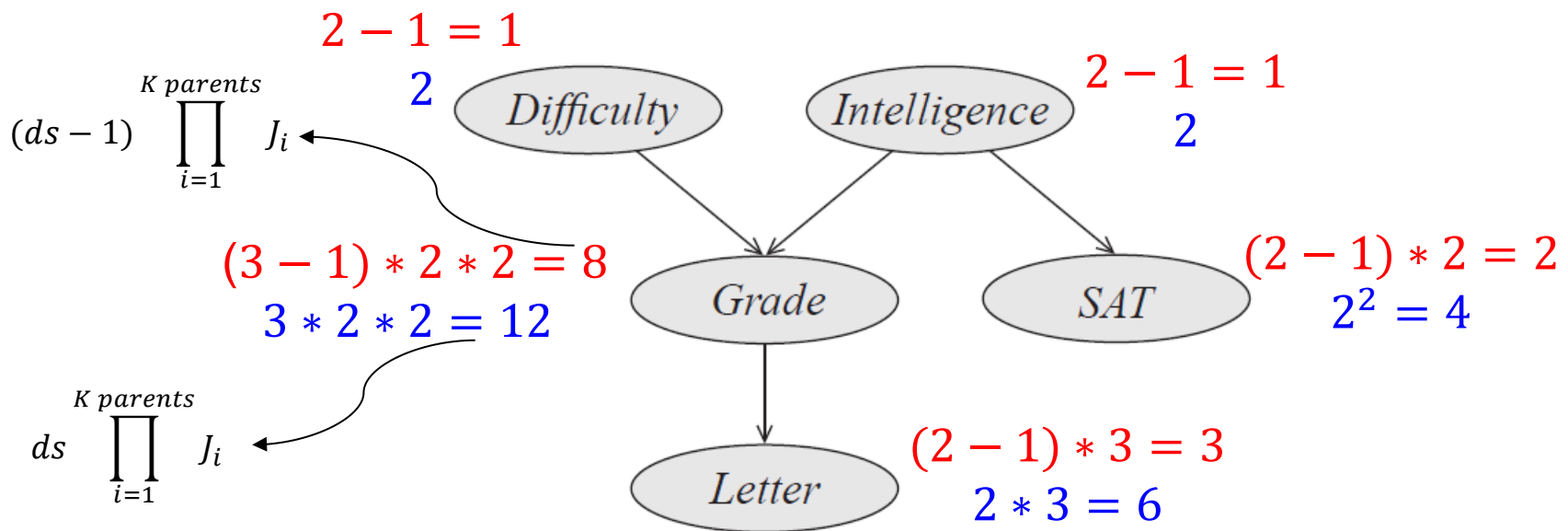
# Number of Probabilities in Bayes Nets

The diagram below should confirm the calculations in the previous slide.



# Number of Probabilities in Bayes Nets

**[Question]** How many probabilities are required by all CPTs of the Bayesian Network below considering that all variables except  $G$  are binary— $G$ 's domain size is 3?



Notation:

$ds$  = domain size of variable with  $K$  parents.

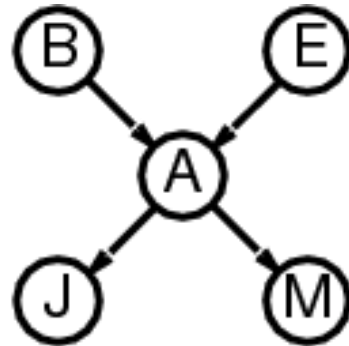
$J_i$  = domain size of parent rand. variable  $i$

Concise version:  $1 + 1 + 8 + 2 + 3 = 15$

Full enumeration:  $2 + 2 + 12 + 4 + 6 = 26$

# Global Semantics

- “Global” semantics refers to the full joint distribution as the product of local conditional distributions:



- $P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$
- Example:  $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e) =$   
 $P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) = 0.9 \times$   
 $0.7 \times 0.001 \times 0.999 \times 0.998 \approx 0.00063$

# Parameter Learning via MLE (Maximum Likelihood Estimation)

For Conditional Probability Tables (CPTs) with one variable we use  $P(X = x) = \frac{\text{count}(x)+1}{\text{count}(X)+|X|}$ , where  $|X|$ =domain size of variable X

play	P(play)
yes	$(9+1)/(14+2)=0.625$
no	$(5+1)/(14+2)=0.375$

For CPTs with two variables we use  $P(x|y) = \frac{\text{count}(x|y)+1}{\text{count}(y)+|X|}$

outlook	play	P(outlook)
sunny	yes	$(2+1)/(9+3)=0.25$
overcast	yes	$(4+1)/(9+3)=0.417$
rainy	yes	$(3+1)/(9+3)=0.333$
sunny	no	$(3+1)/(5+3)=0.5$
overcast	no	$(0+1)/(5+3)=0.125$
rainy	no	$(2+1)/(5+3)=0.375$

For CPTs with 3 vars. we use  $P(x|y, z) = \frac{\text{count}(x|y, z)+1}{\text{count}(y, z)+|X|}$ , and so on

# Techniques for Parameter Learning avoiding Zero Probabilities

## 1. Laplace smoothing

$P(x) = \frac{\text{count}(x)+1}{N+J}$ , where  $N$  is the total number of data points and  $J$  is the total number of possible outcomes (domain size).

## 2. Additive smoothing

$P(x) = \frac{\text{count}(x)+l}{N+l*J}$ , where  $0 < l < 1$ .

## 3. Dirichlet priors

A Dirichlet prior is a probability distribution over the parameters of a discrete distribution. The prior ensures that all events have non-zero probabilities by distributing probability mass across all possible events.

# Techniques for Parameter Learning avoiding Zero Probabilities

## 1. Laplace smoothing

$P(x) = \frac{\text{count}(x)+1}{N+J}$ , where  $N$  is the total number of data points and  $J$  is the total number of possible outcomes (domain size).

## 2. Additive smoothing

$P(x) = \frac{\text{count}(x)+l}{N+l*J}$ , where  $0 < l < 1$ .

Look for an implementation of MLE with Laplace/Additive smoothing during this week's workshop:  
**CPT\_Generator.py**

## 3. Dirichlet priors (example in appendix 2)

A Dirichlet prior is a probability distribution over the parameters of a discrete distribution. The prior ensures that all events have non-zero probabilities by distributing probability mass across all possible events.

# Techniques for Parameter Learning with Missing Data

- Remove data with missing values
- Probabilistic inference, to predict missing values:
  - Step 1: train a model on observable data.
  - Step 2: use model from step 1 to fill missing values.
  - Step 3: train a new model on fully labelled data.
- Using the EM algorithm (can be slow)
  - Initialisation (random or MLE on observable data)
  - Expectation step: estimate missing/expected counts.
  - Maximisation step: MLE using expected counts.
  - Iterate E and M steps until convergence

EM Reference:  
Koller & Friedman  
2009. [Section 19.2.2](#)

# Today

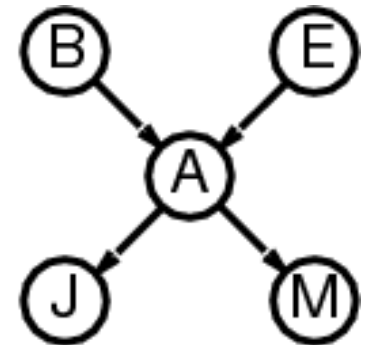
- Introduction to Bayesian networks
  - Graphical representation
  - Probabilistic representation
  - Parameter learning
- **Algorithms for exact inference**
  - Inference by enumeration
  - Inference by variable elimination



# Inference by Enumeration

- Sums out variables from the joint without actually constructing its explicit representation.

- Simple query on the burglary network:

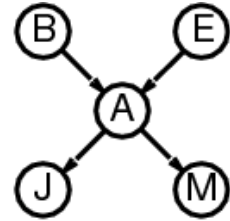


- $P(B|j, m) = \frac{P(B, j, m)}{P(j, m)}$
- $P(B|j, m) = \alpha P(B, j, m)$
- $P(B|j, m) = \alpha \sum_a \sum_e P(B, e, a, j, m)$

└──────────> Normalisation constant

# Inference by Enumeration

$$P(B|j, m) = \alpha \sum_a \sum_e P(B, e, a, j, m)$$



Rewriting joint entries using product of CPT entries:

$$P(B|j, m) = \alpha \sum_a \sum_e P(B)P(e)P(a|b, e)P(j|a)P(m|a)$$

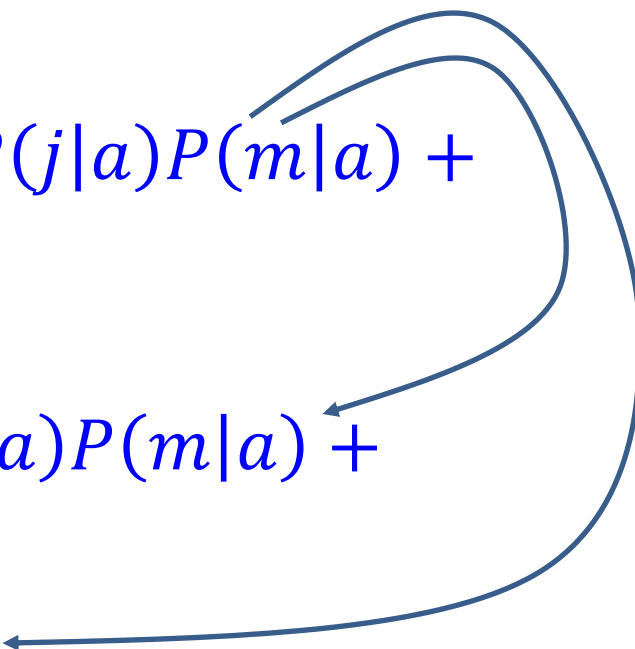
$$= \alpha P(B) \sum_e P(e) \sum_a P(a|b, e)P(j|a)P(m|a)$$

$$= \alpha < P(b|j, m), P(\neg b|j, m) >$$

# Inference by Enumeration: $P(b|j, m)$

$$\begin{aligned} P(b|j, m) &= \alpha \sum_a \sum_e P(b)P(e)P(a|b, e)P(j|a)P(m|a) \\ &= \alpha P(b) \sum_e P(e) \sum_a P(a|b, e)P(j|a)P(m|a) \end{aligned}$$

$$= \alpha P(b) \sum_e P(e) [P(a|b, e)P(j|a)P(m|a) + P(\neg a|b, e)P(j|\neg a)P(m|\neg a)]$$

$$\begin{aligned} &= \alpha P(b) [\textcolor{red}{P(e)} [P(a|b, e)P(j|a)P(m|a) + \\ &P(\neg a|b, e)P(j|\neg a)P(m|\neg a)] + \\ &\textcolor{red}{P(\neg e)} [P(a|b, \neg e)P(j|a)P(m|a) + \\ &P(\neg a|b, \neg e)P(j|\neg a)P(m|\neg a)]] \end{aligned}$$


# Inference by Enumeration: $P(b|j, m)$

$$\begin{aligned} P(b|j, m) = & \alpha P(b) [ \textcolor{red}{P(e)} [P(a|b, e)P(j|a)P(m|a) + \\ & P(\neg a|b, e)P(j|\neg a)P(m|\neg a)] + \\ & \textcolor{red}{P(\neg e)} [P(a|b, \neg e)P(j|a)P(m|a) + \\ & P(\neg a|b, \neg e)P(j|\neg a)P(m|\neg a)] ] \end{aligned}$$

$$\begin{aligned} = & \alpha [ 0.001 \times [ 0.002 \times [0.95 \times 0.9 \times 0.7 + 0.05 \times 0.05 \\ & \times 0.01] + [0.998 \times [0.94 \times 0.9 \times 0.7 + 0.06 \times 0.05 \\ & \times 0.01]] ] \end{aligned}$$

$$\begin{aligned} = & \alpha [0.001 \times [ 0.002 \times [0.5985 + 0.000025] + 0.998 \\ & \times [0.5922 + 0.000003]] ] \end{aligned}$$

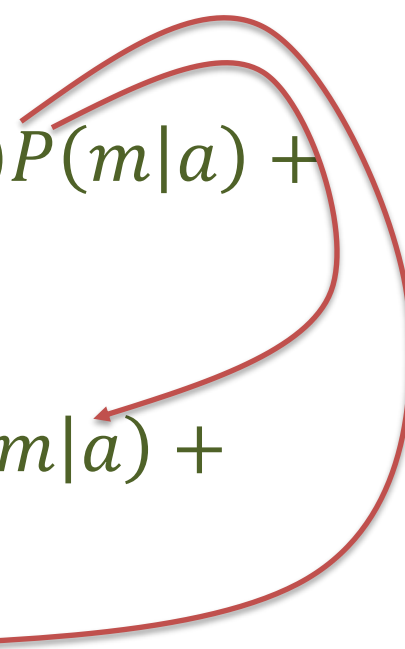
$$= \alpha [0.001 \times [ 0.0001197 + 0.591045]]$$

$$= \alpha 0.000592243$$

# Inference by Enumeration: $P(\neg b|j, m)$

$$\begin{aligned} P(\neg b|j, m) &= \alpha \sum_a \sum_e P(\neg b)P(e)P(a|\neg b, e)P(j|a)P(m|a) \\ &= \alpha P(\neg b) \sum_e P(e) \sum_a P(a|\neg b, e)P(j|a)P(m|a) \end{aligned}$$

$$= \alpha P(\neg b) \sum_e P(e) [P(a|\neg b, e)P(j|a)P(m|a) + P(\neg a|\neg b, e)P(j|\neg a)P(m|\neg a)]$$

$$\begin{aligned} &= \alpha P(\neg b) [P(e) [P(a|\neg b, e)P(j|a)P(m|a) + \\ &P(\neg a|\neg b, e)P(j|\neg a)P(m|\neg a)] + \\ &P(\neg e) [P(a|\neg b, \neg e)P(j|a)P(m|a) + \\ &P(\neg a|\neg b, \neg e)P(j|\neg a)P(m|\neg a)]] \end{aligned}$$


# Inference by Enumeration: $P(\neg b|j, m)$

$$\begin{aligned} P(\neg b|j, m) = & \alpha P(\neg b) [P(e) [P(a|\neg b, e)P(j|a)P(m|a) + \\ & P(\neg a|\neg b, e)P(j|\neg a)P(m|\neg a)] + \\ & P(\neg e) [P(a|\neg b, \neg e)P(j|a)P(m|a) + \\ & P(\neg a|\neg b, \neg e)P(j|\neg a)P(m|\neg a)]] \end{aligned}$$

$$\begin{aligned} = & \alpha [ 0.999 \times [ 0.002 \times [0.29 \times 0.9 \times 0.7 + 0.71 \times 0.05 \\ & \times 0.01] + [0.998 \times [0.001 \times 0.9 \times 0.7 + 0.999 \times 0.05 \\ & \times 0.01]]] \end{aligned}$$

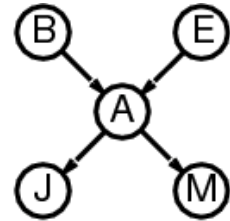
$$\begin{aligned} = & \alpha [0.999 \times [ 0.002 \times [0.1827 + 0.000355] + 0.998 \\ & \times [0.00063 + 0.0004995]]] \end{aligned}$$

$$= \alpha [0.999 \times [0.00036611 + 0.00112724]]$$

$$= \alpha 0.001491858$$

# Inference by Enumeration: $P(B|j, m)$

$$P(B|j, m) = \alpha \sum_a \sum_e P(B, e, a, j, m)$$



Rewriting joint entries using product of CPT entries:

$$P(B|j, m) = \alpha \sum_a \sum_e P(B)P(e)P(a|b, e)P(j|a)P(m|a)$$

$$= \alpha P(B) \sum_e P(e) \sum_a P(a|b, e)P(j|a)P(m|a)$$

$$= \alpha \langle P(b|j, m), P(\neg b|j, m) \rangle$$

$$= \alpha \langle 0.000592243, 0.001491858 \rangle$$

$$= \langle 0.2842, 0.7158 \rangle$$

$$\alpha = \frac{1}{0.000592243 + 0.001491858} = 479.82$$

# Inference by Enumeration: *Algorithm*

**function** ENUMERATION-ASK( $X, \mathbf{e}, bn$ ) **returns** a distribution over  $X$

**inputs:**  $X$ , the query variable

$\mathbf{e}$ , observed values for variables  $\mathbf{E}$

$bn$ , a Bayesian network with variables  $\{X\} \cup \mathbf{E} \cup \mathbf{Y}$

$Q(X) \leftarrow$  a distribution over  $X$ , initially empty

**for each** value  $x_i$  of  $X$  **do**

    extend  $\mathbf{e}$  with value  $x_i$  for  $X$

$Q(x_i) \leftarrow$  ENUMERATE-ALL(VARS[ $bn$ ],  $\mathbf{e}$ )

**return** NORMALIZE( $Q(X)$ )

---

**function** ENUMERATE-ALL( $vars, \mathbf{e}$ ) **returns** a real number

**if** EMPTY?( $vars$ ) **then return** 1.0

$Y \leftarrow$  FIRST( $vars$ )

**if**  $Y$  has value  $y$  in  $\mathbf{e}$

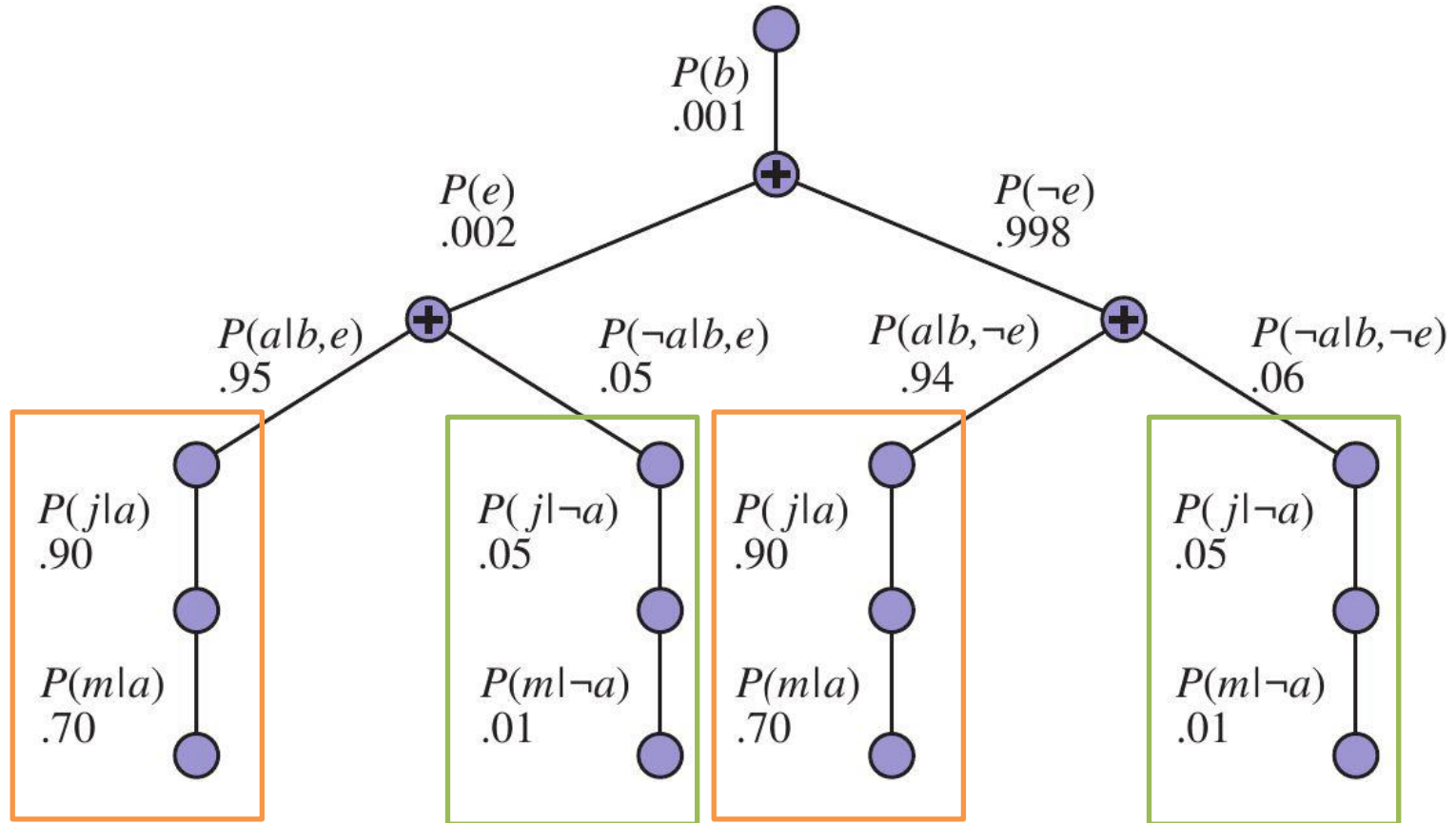
**then return**  $P(y \mid Pa(Y)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}$ )

**else return**  $\sum_y P(y \mid Pa(Y)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}_y$ )

        where  $\mathbf{e}_y$  is  $\mathbf{e}$  extended with  $Y = y$



# Inference by Enumeration: Evaluation Tree



Enumeration can be inefficient due to repeated computations

# Today

- Introduction to Bayesian networks
  - Graphical representation
  - Probabilistic representation
  - Parameter learning
- Algorithms for exact inference
  - Inference by enumeration
  - **Inference by variable elimination (appendix 1)**

# Homework (recommended)

1. Calculate  $P(E | j, m)$  using **inference by enumeration** with pen and paper—revising the example provided above.
2. Calculate  $P(E | j, m)$  using **variable elimination** with pen and paper—but first look at the example in appendix 1.

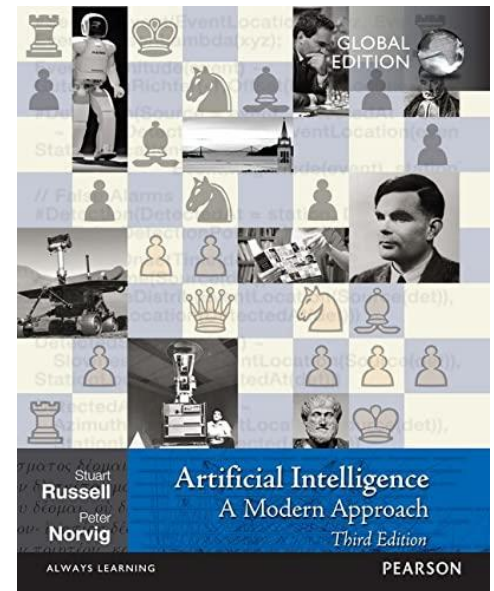
# Today

- Introduction to Bayesian networks
- Parameter learning via Maximum Likelihood Estimation (MLE) with smoothing techniques
- Inference by enumeration
- Inference by variable elimination

## Readings:

Russell & Norvig 2016. [Chapters 14-14.4](#)

Koller & Friedman 2009. Section 17.3.2



# This and Next Week

## **Workshop (tomorrow):**

Exercises using Bayesian networks  
Python program for exact inference

## **Lecture (next week):**

Structure Learning for Bayesian Networks

Reading: [Kitson et al. A survey on Bayesian Network structure learning, 2023](#)

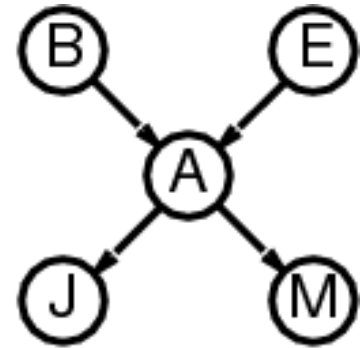
Questions?

# Appendix 1

## Probabilistic Inference via Variable Elimination

# Inference by Variable Elimination

- Idea:
  - Do the calculation once; and
  - Save the results for later use.



- Variable elimination evaluates expressions in right-to-left order, and uses factors  $f_i$  (matrices) as follows:

$$P(B|j, m) = \alpha \underbrace{P(B)}_{f_1(B)} \underbrace{\sum_e P(e)}_{f_2(E)} \underbrace{\sum_a P(a|b, e)}_{f_3(A, B, E)} \underbrace{P(j|a)}_{f_4(A)} \underbrace{P(m|a)}_{f_5(A)}$$

# Inference by Variable Elimination

$$f_4(A) = \langle P(j|a), P(j|\neg a) \rangle = \langle 0.90, 0.05 \rangle$$

$$f_5(A) = \langle P(m|a), P(m|\neg a) \rangle = \langle 0.70, 0.01 \rangle$$

Therefore,  $P(B|j, m) =$

$$\alpha f_1(B) \times \sum_e f_2(E) \times \sum_a \underbrace{f_3(A, B, E) \times f_4(A) \times f_5(A)}_{f_6(B, E)},$$

where  $\times$  denotes a pointwise product operation.

$$\begin{aligned} f_6(B, E) &= \sum_a f_3(A, B, E) \times f_4(A) \times f_5(A) \\ &= [f_3(a, B, E) \times f_4(a) \times f_5(a)] + [f_3(\neg a, B, E) \times f_4(\neg a) \times f_5(\neg a)] \end{aligned}$$



# Inference by Variable Elimination

$$\text{Therefore, } P(B|j, m) = \alpha f_1(B) \times \underbrace{\sum_e f_2(E) \times f_6(B, E)}_{f_7(B)}$$

Summing out  $E$  we get:

$$\begin{aligned} f_7(B) &= \sum_e f_2(E) \times f_6(B, E) \\ &= [f_2(e) \times f_6(b, e)] + [f_2(\neg e) \times f_6(b, \neg e)] \end{aligned}$$

$$\text{Thus, } P(B|j, m) = \alpha f_1(B) \times f_7(B)$$

We only need to know how to do operations with factors!

# Pointwise Product with Factors

$A$	$B$	$\mathbf{f}_1(A, B)$	$B$	$C$	$\mathbf{f}_2(B, C)$	$A$	$B$	$C$	$\mathbf{f}_3(A, B, C)$
T	T	.3	T	T	.2	T	T	T	$.3 \times .2 = .06$
T	F	.7	T	F	.8	T	T	F	$.3 \times .8 = .24$
F	T	.9	F	T	.6	T	F	T	$.7 \times .6 = .42$
F	F	.1	F	F	.4	T	F	F	$.7 \times .4 = .28$
						F	T	T	$.9 \times .2 = .18$
						F	T	F	$.9 \times .8 = .72$
						F	F	T	$.1 \times .6 = .06$
						F	F	F	$.1 \times .4 = .04$

**Figure 14.10** Illustrating pointwise multiplication:  $\mathbf{f}_1(A, B) \times \mathbf{f}_2(B, C) = \mathbf{f}_3(A, B, C)$ .

# Operations on Factors

$A$	$B$	$\mathbf{f}_1(A, B)$	$B$	$C$	$\mathbf{f}_2(B, C)$	$A$	$B$	$C$	$\mathbf{f}_3(A, B, C)$
T	T	.3	T	T	.2	T	T	T	$.3 \times .2 = .06$
T	F	.7	T	F	.8	T	T	F	$.3 \times .8 = .24$
F	T	.9	F	T	.6	T	F	T	$.7 \times .6 = .42$
F	F	.1	F	F	.4	T	F	F	$.7 \times .4 = .28$
						F	T	T	$.9 \times .2 = .18$
						F	T	F	$.9 \times .8 = .72$
						F	F	T	$.1 \times .6 = .06$
						F	F	F	$.1 \times .4 = .04$

**Figure 14.10** Illustrating pointwise multiplication:  $\mathbf{f}_1(A, B) \times \mathbf{f}_2(B, C) = \mathbf{f}_3(A, B, C)$ .

$$\begin{aligned}
 f(Y, Z) &= \sum_x f(X, Y, Z) = f(x, Y, Z) + f(\neg x, Y, Z) \\
 &= \begin{pmatrix} 0.06 & 0.24 \\ 0.42 & 0.28 \end{pmatrix} + \begin{pmatrix} 0.18 & 0.72 \\ 0.06 & 0.04 \end{pmatrix} = \begin{pmatrix} 0.24 & 0.96 \\ 0.48 & 0.32 \end{pmatrix}
 \end{aligned}$$

# Inference by Variable Elimination: Full Example

$$\begin{aligned}
 P(B|j, m) &= \alpha \underbrace{P(B)}_{f_1(B)} \underbrace{\sum_e P(e)}_{f_2(E)} \underbrace{\sum_a P(a|b, e)}_{f_3(A, B, E)} \underbrace{P(j|a)}_{f_4(A)} \underbrace{P(m|a)}_{f_5(A)} \\
 &= \alpha f_1(B) \times \sum_e f_2(E) \times \underbrace{\sum_a f_3(A, B, E) \times f_4(A) \times f_5(A)}_{f_6(B, E)}
 \end{aligned}$$

$$\begin{aligned}
 f_6(B, E) &= \\
 &= [f_3(a, B, E) \times f_4(a) \times f_5(a)] + [f_3(\neg a, B, E) \times f_4(\neg a) \times f_5(\neg a)] \\
 &= \begin{pmatrix} B & E & f_3 \\ t & t & 0.95 \\ t & f & 0.94 \\ f & t & 0.29 \\ f & f & 0.001 \end{pmatrix} \times 0.63 + \begin{pmatrix} B & E & f_4 \\ t & t & 0.05 \\ t & f & 0.06 \\ f & t & 0.71 \\ f & f & 0.94 \end{pmatrix} \times 0.0005 = \begin{pmatrix} B & E & f_6 \\ t & t & 0.59852 \\ t & f & 0.59222 \\ f & t & 0.18305 \\ f & f & 0.00110 \end{pmatrix}
 \end{aligned}$$

# Inference by Variable Elimination: Full Example

$$\begin{aligned}
 f_7(B) &= [f_2(e)f_6(B, e)] + [f_2(\neg e)f_6(B, \neg e)] \\
 &= 0.002 \times \begin{pmatrix} B & f_6 \\ t & 0.59852 \\ f & 0.18305 \end{pmatrix} + 0.998 \begin{pmatrix} B & f_6 \\ t & 0.59222 \\ f & 0.00110 \end{pmatrix} = \begin{pmatrix} B & f_7 \\ t & 0.59223 \\ f & 0.00146 \end{pmatrix}
 \end{aligned}$$

$$\begin{aligned}
 P(B|j, m) &= \alpha f_1(B) \times f_7(B) \\
 &= \alpha \begin{pmatrix} B & f_1 \\ t & 0.001 \\ f & 0.999 \end{pmatrix} \times \begin{pmatrix} B & f_7 \\ t & 0.59223 \\ f & 0.00146 \end{pmatrix} = \alpha \begin{pmatrix} P(B|j, m) \\ t & 0.000592 \\ f & 0.001458 \end{pmatrix} \\
 &= \langle 0.289, 0.711 \rangle
 \end{aligned}$$

$\alpha = \frac{1}{0.000592 + 0.001458}$

# Variable Elimination: *Algorithm*

**function** ELIMINATION-ASK( $X, \mathbf{e}, bn$ ) **returns** a distribution over  $X$

**inputs:**  $X$ , the query variable

$\mathbf{e}$ , observed values for variables  $\mathbf{E}$

$bn$ , a Bayesian network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$

$factors \leftarrow []$

**for each**  $var$  **in** ORDER( $bn.VARS$ ) **do**

$factors \leftarrow [\text{MAKE-FACTOR}(var, \mathbf{e}) | factors]$

**if**  $var$  is a hidden variable **then**  $factors \leftarrow \text{SUM-OUT}(var, factors)$

**return** NORMALIZE(PPOINTWISE-PRODUCT( $factors$ ))

# Appendix 2

## Maximum Likelihood Estimation (MLE) with Dirichlet Priors

# Dirichlet Priors via Moment Matching

- Compute empirical probabilities and variances

$\hat{p}_i = \frac{\text{count}(x_i)}{N}$ , where  $N$ =the total number of data points of interest.

$\hat{\sigma}_i^2 = \frac{\hat{p}_i(1-\hat{p}_i)}{N}$ , which is the empirical variance of probability  $\hat{p}_i$ .

- Match the moments

Mean:  $E[P(X = x_i)] = \frac{\alpha_i}{\sum_j \alpha_j}$

Variance:  $\text{Var}(P(X = x_i)) = \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}$ , where  $\alpha_0 = \sum_j \alpha_j$



# Estimation of Dirichlet Parameters: PlayTennis and Outlook Data (1/4)

Compute empirical probabilities:

$$\hat{p}_{yes,sunny} = \frac{\text{count}(yes, sunny)}{\text{count}(yes)} = \frac{2}{9} = 0.222$$

$$\hat{p}_{yes,overcast} = \frac{\text{count}(yes, overcast)}{\text{count}(yes)} = \frac{4}{9} = 0.444$$

$$\hat{p}_{yes,rain} = \frac{\text{count}(yes, rain)}{\text{count}(yes)} = \frac{3}{9} = 0.333$$

$$\hat{p}_{no,sunny} = \frac{\text{count}(no, sunny)}{\text{count}(no)} = \frac{3}{5} = 0.6$$

$$\hat{p}_{no,overcast} = \frac{\text{count}(no, overcast)}{\text{count}(no)} = \frac{0}{5} = 0$$

$$\hat{p}_{no,rain} = \frac{\text{count}(no, rain)}{\text{count}(no)} = \frac{2}{5} = 0.4$$

# Estimation of Dirichlet Parameters: PlayTennis and Outlook Data (2/4)

Compute empirical variances:

$$\hat{\sigma}_{yes,sunny}^2 = \frac{\hat{p}_{yes,sunny}(1 - \hat{p}_{yes,sunny})}{count(yes)} = \frac{0.222 * 0.778}{9} = 0.0192$$

$$\hat{\sigma}_{yes,overcast}^2 = \frac{\hat{p}_{yes,overcast}(1 - \hat{p}_{yes,overcast})}{count(yes)} = \frac{0.444 * 0.556}{9} = 0.0274$$

$$\hat{\sigma}_{yes,rain}^2 = \frac{\hat{p}_{yes,rain}(1 - \hat{p}_{yes,rain})}{count(yes)} = \frac{0.333 * 0.667}{9} = 0.0247$$

$$\hat{\sigma}_{no,sunny}^2 = \frac{\hat{p}_{no,sunny}(1 - \hat{p}_{no,sunny})}{count(no)} = \frac{0.6 * 0.4}{5} = 0.048$$

$$\hat{\sigma}_{no,overcast}^2 = \frac{\hat{p}_{no,overcast}(1 - \hat{p}_{no,overcast})}{count(no)} = \frac{0 * 1}{5} = 0$$

$$\hat{\sigma}_{no,rain}^2 = \frac{\hat{p}_{no,rain}(1 - \hat{p}_{no,rain})}{count(no)} = \frac{0.4 * 0.6}{5} = 0.048$$

# Estimation of Dirichlet Parameters: PlayTennis and Outlook Data (3/4)

- From moment matching we know that

$$\hat{p}_i = \frac{\alpha_i}{\alpha_0} \text{ and that } \hat{\sigma}_i^2 = \frac{\hat{p}_i(1-\hat{p}_i)}{\alpha_0+1} \Rightarrow \alpha_0 = \frac{\hat{p}_i(1-\hat{p}_i)}{\hat{\sigma}_i^2} - 1$$

- Estimating  $\alpha_0$  for *PlayTennis* = *yes*:

$$\alpha_0^{yes} = \frac{\hat{p}_{yes,sunny}(1 - \hat{p}_{yes,sunny})}{\hat{\sigma}_{yes,sunny}^2} - 1 = \frac{0.222 * 0.778}{0.0192} - 1 = 8$$

- Estimating  $\alpha_0$  for *PlayTennis* = *no*:

$$\alpha_0^{no} = \frac{\hat{p}_{no,sunny}(1 - \hat{p}_{no,sunny})}{\hat{\sigma}_{no,sunny}^2} - 1 = \frac{0.4 * 0.6}{0.048} - 1 = 4$$

# Estimation of Dirichlet Parameters: PlayTennis and Outlook Data (4/4)

- Dirichlet parameters  $\alpha_i$  for  $PlayTennis = yes$  :

$$\alpha_{yes,sunny} = \hat{p}_{yes,sunny} * \alpha_0^{yes} = 0.222 * 8 = 1.78$$

$$\alpha_{yes,overcast} = \hat{p}_{yes,overcast} * \alpha_0^{yes} = 0.444 * 8 = 3.55$$

$$\alpha_{yes,rain} = \hat{p}_{yes,rain} * \alpha_0^{yes} = 0.333 * 8 = 2.66$$

- Dirichlet parameters  $\alpha_i$  for  $PlayTennis = no$  :

$$\alpha_{no,sunny} = \hat{p}_{no,sunny} * \alpha_0^{no} = 0.6 * 4 = 2.4$$

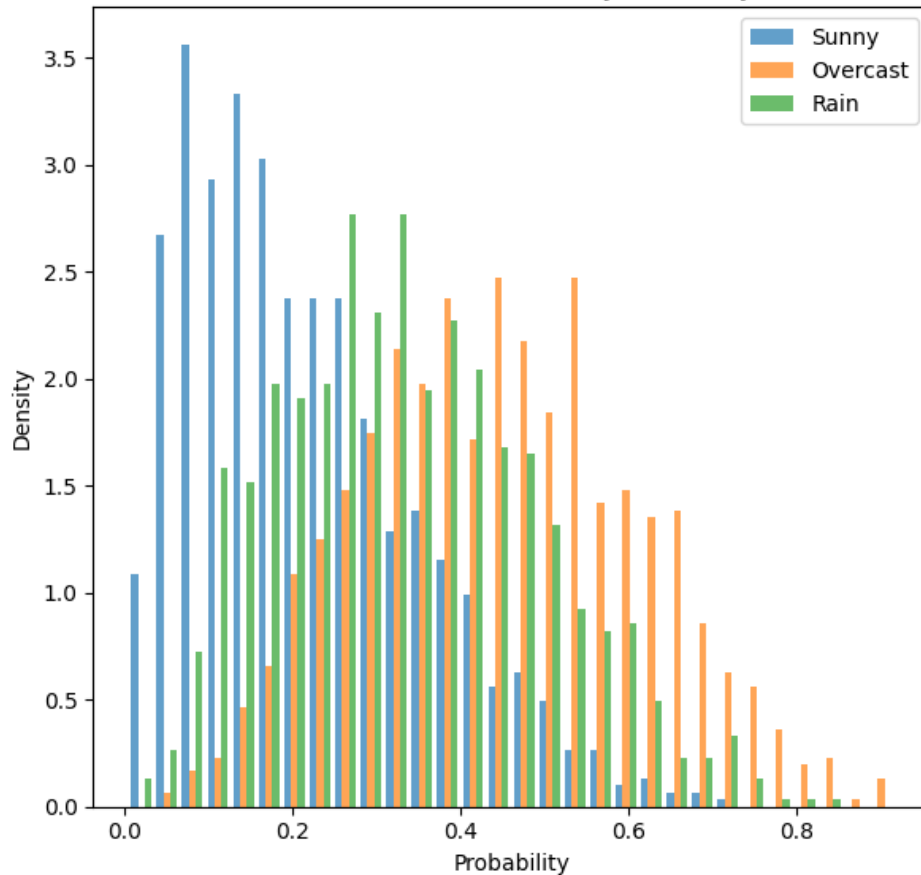
$$\alpha_{no,overcast} = \hat{p}_{no,overcast} * \alpha_0^{no} = 0 * 4 + \epsilon = 0.5$$

$$\alpha_{no,rain} = \hat{p}_{no,rain} * \alpha_0^{no} = 0.4 * 4 = 1.6$$

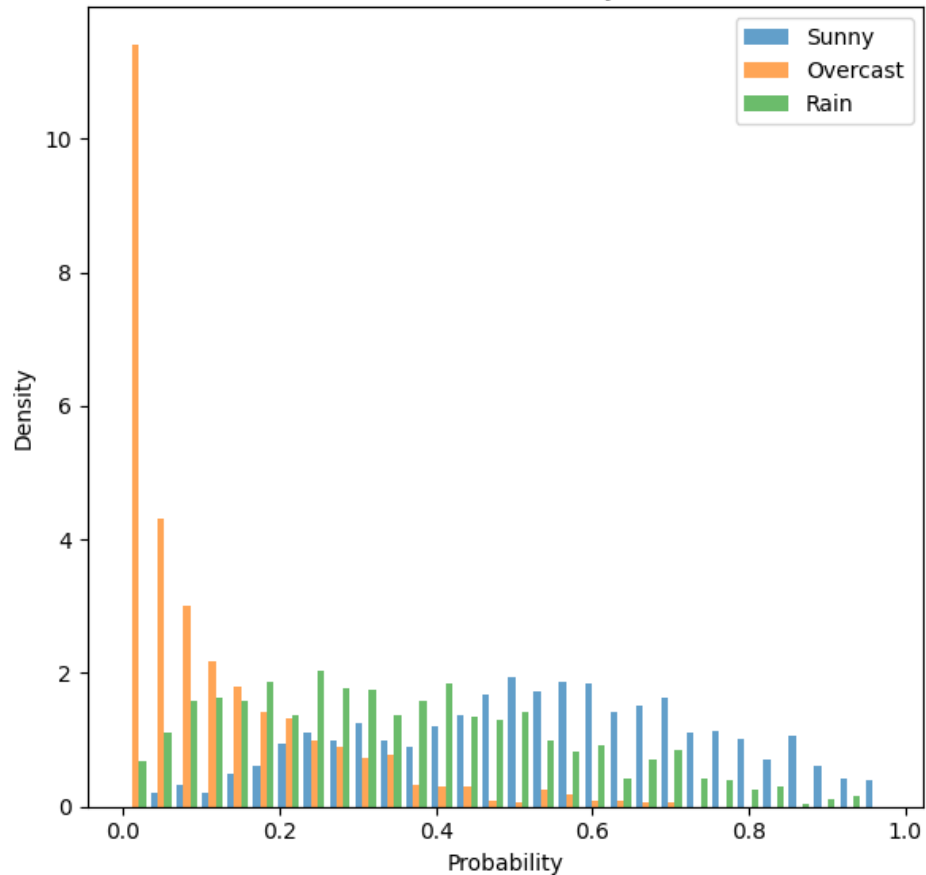
where  $\epsilon = 0.5$  is used to avoid zero values. While higher values of  $\alpha$  mean more confidence in the estimated probabilities, lower values suggest less confidence or more uncertainty in the probabilities.

# Dirichlet Distributions for PlayTennis and Outlook Example (1K samples)

Dirichlet Distribution for PlayTennis = yes



Dirichlet Distribution for PlayTennis = no



```
samples_yes = np.random.dirichlet([1.78, 3.55, 2.66], 1000) # 30 bins  
samples_no = np.random.dirichlet([2.4, 0.5, 1.6], 1000) # 30 bins
```

# MLE with Dirichlet Parameters

$$P(\text{sunny}|\text{yes}) = \frac{\text{count}(\text{yes}, \text{sunny}) + \alpha_{\text{yes}, \text{sunny}}}{\text{count}(\text{yes}) + \alpha_0(\text{yes})} = \frac{2 + 1.78}{9 + 8} = 0.2223$$

$$P(\text{overcast}|\text{yes}) = \frac{\text{count}(\text{yes}, \text{overcast}) + \alpha_{\text{yes}, \text{overcast}}}{\text{count}(\text{yes}) + \alpha_0(\text{yes})} = \frac{4 + 3.55}{9 + 8} = 0.4441$$

$$P(\text{rain}|\text{yes}) = \frac{\text{count}(\text{yes}, \text{rain}) + \alpha_{\text{yes}, \text{rain}}}{\text{count}(\text{yes}) + \alpha_0(\text{yes})} = \frac{3 + 2.66}{9 + 8} = 0.3329$$

$$P(\text{sunny}|\text{no}) = \frac{\text{count}(\text{no}, \text{sunny}) + \alpha_{\text{no}, \text{sunny}}}{\text{count}(\text{no}) + \alpha_0(\text{no})} = \frac{3 + 2.4}{5 + 4.5} = 0.5684$$

$$P(\text{overcast}|\text{no}) = \frac{\text{count}(\text{no}, \text{overcast}) + \alpha_{\text{no}, \text{overcast}}}{\text{count}(\text{no}) + \alpha_0(\text{no})} = \frac{0 + 0.5}{5 + 4.5} = 0.0526$$

$$P(\text{rain}|\text{no}) = \frac{\text{count}(\text{no}, \text{rain}) + \alpha_{\text{no}, \text{rain}}}{\text{count}(\text{no}) + \alpha_0(\text{no})} = \frac{2 + 1.6}{5 + 4.5} = 0.3789$$

# Is Dirichlet-Based MLE worth it?

- General formula:

$$P(X = x_i | Pa(X) = pa_j) = \frac{\text{count}(x_i, pa_j) + \alpha_i}{\text{count}(pa_j) + \alpha_0}$$

- This is a more advanced and principled approach of parameter learning than MLE with simple smoothing—because it combines observed data (counts) with prior beliefs ( $\alpha_i$ ).
- It can be useful in the presence of small data.