# Learning

- Parameters
- Structure
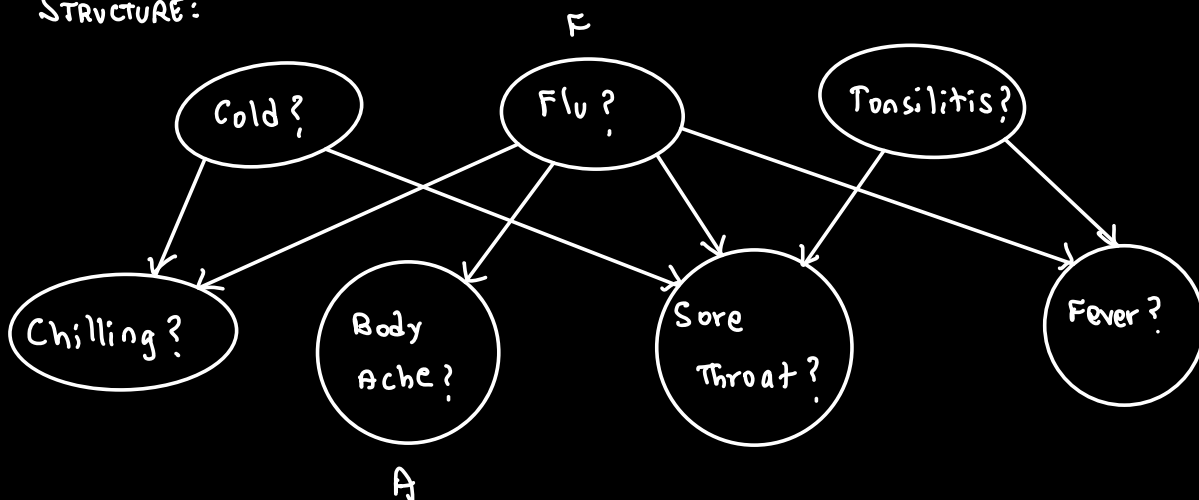
Supervised   vs   Unsupervised
                      ↳ model-oriented
  ↳ Query-oriented


## LEARNING PARAMETERS:

### STRUCTURE:

F

| Cold? | Flu? | Tonsilitis? |

| Chilling? | Body Ache? | Sore Throat? | Fever? |

A

| CPTs can also be estimated from medical records of previous patients | | | | | | |
|---|---|---|---|---|---|---|
| Case | Cold? | Flu? | Tonsillitis? | Chilling? | Bodyache? | Sorethroat? | Fever? |
| 1 | true | false | ? | true | false | false | false → Incomplete example |
| 2 | false | true | false | true | true | false | true |
| 3 | ? | ? | true | false | ? | true | false ↓ |
| : | : | : | : | : | : | : | : Complete example |

examples

Data: Complete, not complete

| F | A | |
|---|---|---|
| t | t | $\theta_{a\mid f}$ |
| t | f | $\theta_{\bar{a}\mid f}$ |
| f | t | $\theta_{a\mid \bar{f}}$ |
| f | f | $\theta_{\bar{a}\mid \bar{f}}$ |

$\nearrow$ Pr $(a\mid f)$

$\left.\begin{array}{c}\\\\\\\end{array}\right\}$ Parameters

$\nwarrow$ Pr $(\bar{a}\mid\bar{f})$

MAXIMUM LIKELIHOOD:

Parameters 1 $\Rightarrow$ BN 1 $\Rightarrow$ $Pr_1(data) = Pr(data\mid parameters_1)$

Parameters 2 $\Rightarrow$ BN2 $\Rightarrow$ $Pr_2(data) = Pr(data\mid parameters_2)$

$e_1 \ e_2 \ ... \ e_n$

LIKELIHOOD $\Rightarrow$ $Pr_1(e_1) \ Pr_1(e_2) ... \ Pr_1(e_n)$     Score 1

$\Rightarrow$ $Pr_2(e_1) \ Pr_2(e_2) \ ... \ Pr_2(e_n)$     Score 2

| Case | H | S | E |
|------|---|---|---|
| 1 | T | F | T |
| 2 | T | F | T |
| 3 | F | T | F |
| 4 | F | F | T |
| 5 | T | F | F |
| 6 | T | F | T |
| 7 | F | F | F |
| 8 | T | F | T |
| 9 | T | F | T |
| 10 | F | F | T |
| 11 | T | F | T |
| 12 | T | T | T |
| 13 | T | F | T |
| 14 | T | T | T |
| 15 | T | F | T |
| 16 | T | F | T |

| H | S | E | $\Pr_D(.)$ |
|---|---|---|-----------|
| T | T | T | 2/16 |
| T | T | F | 0/16 |
| T | F | T | 9/16 |
| T | F | F | 1/16 |
| F | T | T | 0/16 |
| F | T | F | 1/16 |
| F | F | T | 2/16 |
| F | F | F | 1/16 |

Health Aware (H)

Smokes (S)    Exercises (E)

↑                    ↑

DATASET         Empirical

(Complete)      Distribution

| H | |
|---|---|
| t | $\theta_h$ |
| f | $\theta_{\bar{h}}$ |

| H | S | |
|---|---|---|
| t | t | $\theta_{s\mid h}$ |
| t | f | $\theta_{\bar{s}\mid h}$ |
| f | t | $\theta_{s\mid \bar{h}}$ |
| f | f | $\theta_{\bar{s}\mid \bar{h}}$ |

Similarly for H,E

(BAYES CONDITIONING)

$$\theta_{\bar{s}\mid h} = Pr(\bar{s}\mid h) = \frac{Pr(\bar{s}, h)}{Pr(h)}$$

$$= \frac{Pr_D(\bar{s}, h)}{Pr_D(h)}$$
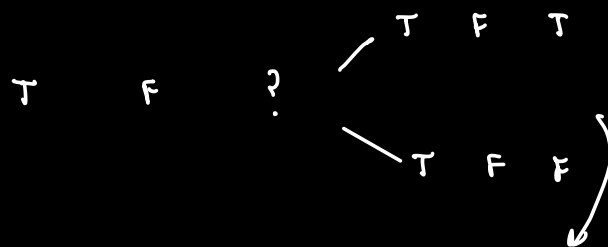
$$= \frac{10/16}{12/16} = 5/6$$

PARAMETER ESTIMATE.

WHAT IF INCOMPLETE DATA?

| H | S | Ĕ |
|---|---|---|
| ⋮ | | |
| T | F | ? |
| ⋮ | | |

EM: [EXPECTATION MAXIMIZATION]

$$CPT_1 \rightsquigarrow BN_1 \rightarrow Pr_1(\cdot)$$

$$T \quad F \quad ?
\begin{cases}
T & F & T \\
T & F & F
\end{cases}$$

$$x = Pr_1(\check{E} = T \mid H = T, S = F)$$

$$y = Pr_1(\check{E} = F \mid H = T, S = F)$$

Fill the $x, y$ in

Recompute CPT

$$CPT_2 \rightarrow BN_2 \rightarrow Pr_2(\cdot)$$

⋮

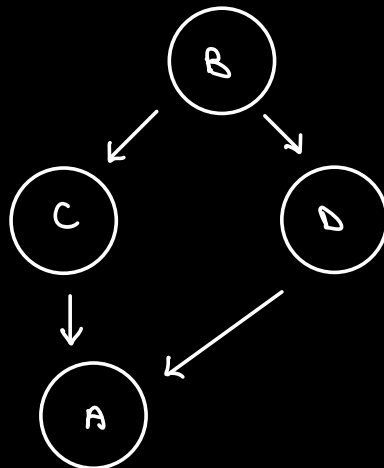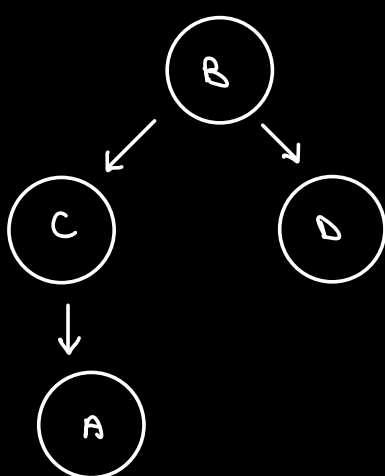$$CPT_3 \rightarrow BN_3 \rightarrow Pr_3(\cdot)$$

⋮

Converges

Likelihood never decreases.

$\quad\quad\quad\quad\quad\quad\quad\quad\hookrightarrow$ Increases or remains same.

Convergence speed $\alpha$ $\dfrac{1}{n_{missing\ data}}$

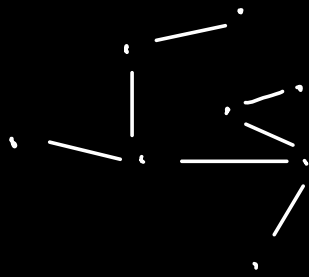So $\quad\left.\begin{array}{l}\theta_1 \\ \\ \theta_2 \\ \\ \theta_3 \\ \quad\vdots \\ \theta_n\end{array}\right\}$ Inference

## LEARNING STRUCTURE:



Choose what gives best score.

# 1. LOCAL SEARCH METHODS:



Approximate

Fast

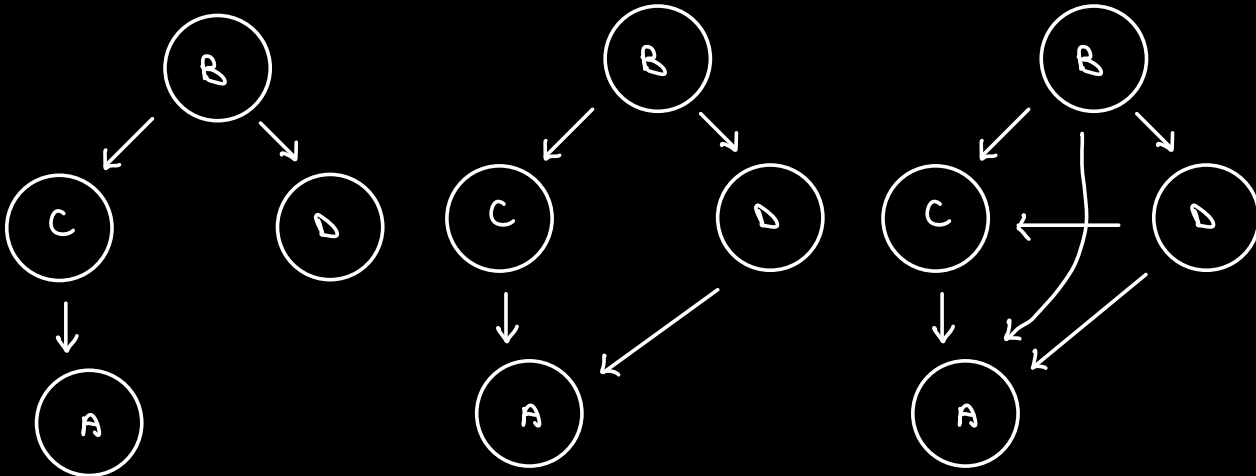Add, remove, reverse an edge

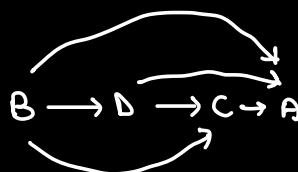# 2. SYSTEMATIC SEARCH METHODS:

A* search

Guaranteed

Slower

# WHY NOT USE MAXIMUM LIKELIHOOD?

Overfitting



always chose

Complete DAG

$$B \longrightarrow D \longrightarrow C \longrightarrow A$$

| x | y |
|---|---|
| 1 | 1.1 |
| 5 | 4.5 |
| 10 | 11 |
| 15 | 14.5 |
| 20 | 22 |

$$y = ax + b$$

$\uparrow \quad \uparrow$

Parameters

to fit perfectly,

we need

$$y = ax^4 + bx^3 + cx^2 + dx + e$$

OVERFITTING

* Not generalizing
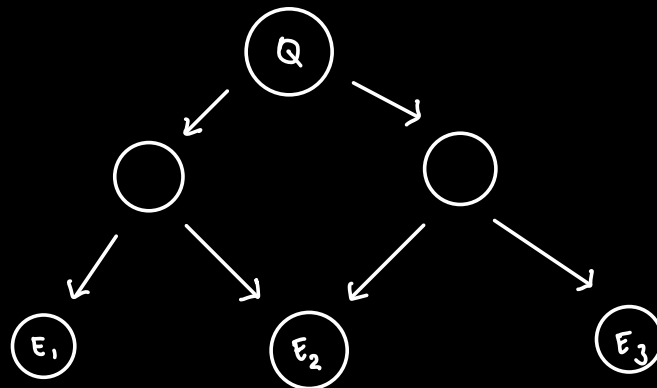
SCORE: Likelihood + Penalty term

$\uparrow$

MDL
Score

$\hookrightarrow$ depends on

$\rightarrow$ number of independent
Parameters

$\rightarrow$ Size of dataset

MODEL - ORIENTED VS QUERY - ORIENTED LEARNING:

Unsupervised vs Supervised Learning

Un-labeled vs labeled data

Model – Oriented:

Learn Structure (can answer any query)

[What we did till now]

Query – Oriented:
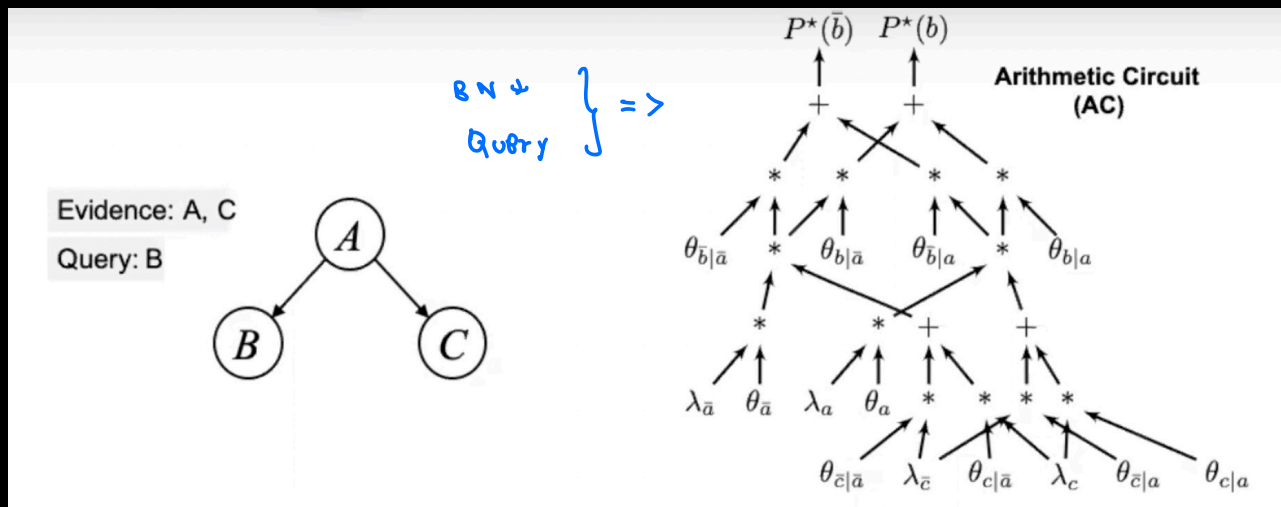
Specific to one query.

$$P(Q | E_1, E_2, E_3)$$

$$\downarrow$$

Query

QUERY – ORIENTED:

| $E_1$ | $E_2$ | $E_3$ | $Q$ |
|---|---|---|---|

↙ labels

BN & Query $\Bigg\}$ =>

$P^\star(\bar{b})$  $P^\star(b)$

**Arithmetic Circuit (AC)**

Evidence: A, C

Query: B

$A$

$B$  $C$

$\theta_{\bar{b}|\bar{a}}$  $\theta_{b|\bar{a}}$  $\theta_{\bar{b}|a}$  $\theta_{b|a}$

$\lambda_{\bar{a}}$  $\theta_{\bar{a}}$  $\lambda_a$  $\theta_a$

$\theta_{\bar{c}|\bar{a}}$  $\lambda_{\bar{c}}$  $\theta_{c|\bar{a}}$  $\lambda_c$  $\theta_{\bar{c}|a}$  $\theta_{c|a}$

Boolean formula
+
weights

$\Downarrow$

NNF CIRCUIT

$\downarrow$

Convert NNF
circuit to AC

$P \to$ Distribution on B (Query)

$\theta \to$ BN Parameters

| A | $\lambda_a$ | $\lambda_{\bar{a}}$ |
|---|---|---|
| T | 1 | 0 |
| F | 0 | 1 |
| ? | 1 | 1 |

Evidence (Input)

A = T , C = F

$\lambda_a = 1$

$\lambda_{\bar{a}} = 0$

$\lambda_c = 0$

$\lambda_{\bar{c}} = 0$

A = F

$\lambda_a = 0$

$\lambda_{\bar{a}} = 1$

$\lambda_c = 1$

$\lambda_{\bar{c}} = 1$
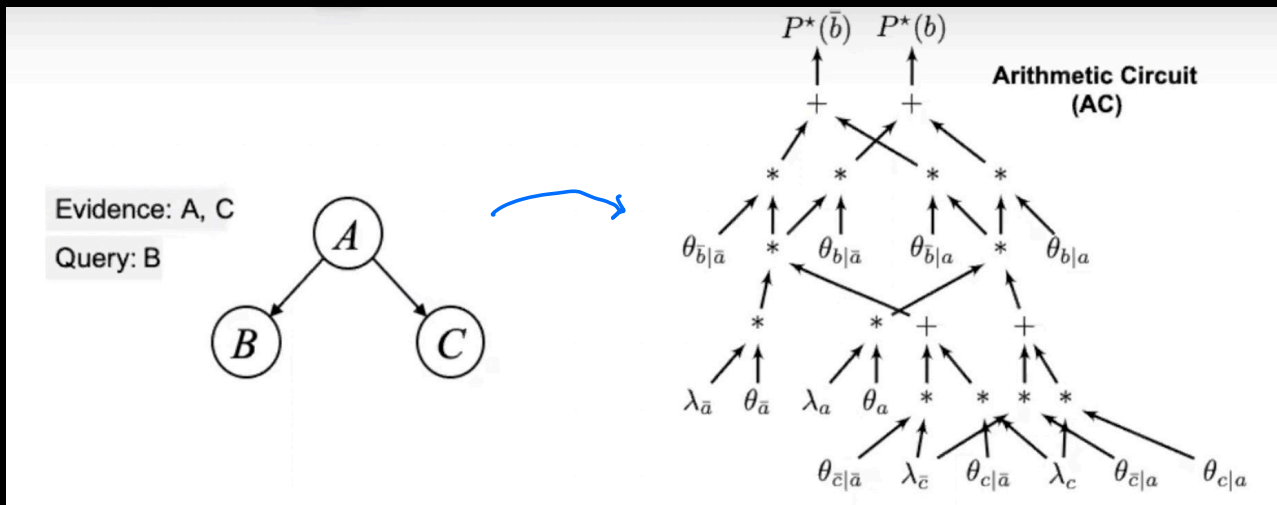
n      variables

d      # values

w      tree width

$$O(n \cdot d^w)$$



Evidence: A, C

Query: B

$P^\star(\bar{b})$  $P^\star(b)$

Arithmetic Circuit (AC)

LABELED DATA

| Input | | Output |
|---|---|---|
| A | C | B |
| T | T | F |
| F | T | F |
| ⋮ | ⋮ | ⋮ |
| T | T | F |

Loss function

(cross entropy)

$P(x)$   $Q(x)$

number (how close they are)

model $\{$    $A = T$ , $C = T$    $\Rightarrow$    $Pr(B)$

                   $\boxed{P}$          $\nwarrow$ Prediction

from
data $\{$    $A = T$ , $C = T$    $\Rightarrow$

| | B | | Label |
|---|---|---|---|
| $\boxed{Q}$ | T | 0 | $\}$ one - hot |
| | F | 1 | distribution |

GRADIENT DESCENT is used to optimize loss function.

     * TensorFlow

     * Pytorch

CROSS ENTROPY :

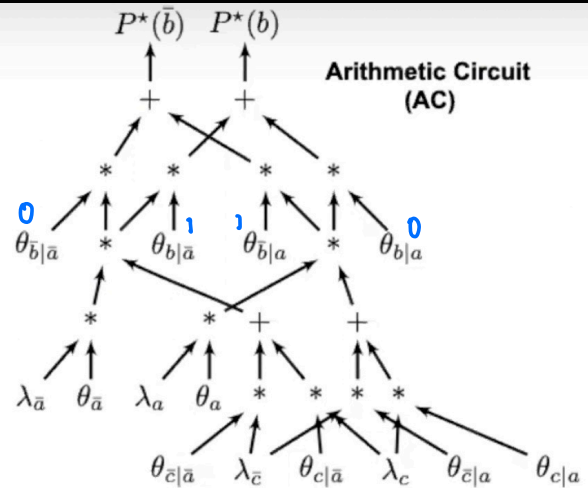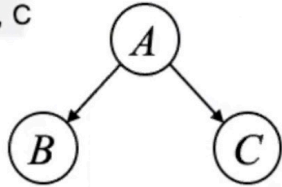     $P(x)$ :   predictions   $\nwarrow$ distribution

     $Q(x)$ :   labels     $\swarrow$

         CE :    $\sum\limits_{x} Q(x) \cdot \log_2 (P(x))$

     LOSS FUNCTION

Evidence: A, C

Query: B

$P^\star(\bar{b})$   $P^\star(b)$

**Arithmetic Circuit (AC)**

$\theta_{\bar{b}|\bar{a}}$   $\theta_{b|\bar{a}}$   $\theta_{\bar{b}|a}$   $\theta_{b|a}$

$\lambda_{\bar{a}}$   $\theta_{\bar{a}}$   $\lambda_a$   $\theta_a$

$\theta_{\bar{c}|\bar{a}}$   $\lambda_{\bar{c}}$   $\theta_{c|\bar{a}}$   $\lambda_c$   $\theta_{\bar{c}|a}$   $\theta_{c|a}$

$B = T$   iff   $A = F$

| A | B | |
|---|---|---|
| t | t | 0 |
| t | f | 1 |
| f | t | 1 |
| f | f | 0 |

Background Knowledge

# FIND THE LARGEST RECTANGLE:



## Rectangle:

Upper left: row, col

Height

Width

Label: Tall or wide



EVIDENCE:

$row_1, row_2 \cdots row_n$ : True, False

$col_1, col_2 \cdots col_n$ : True, False

$pixel_{1,1}$
   row, col

Pr (Tall)     Pr (Wide)    | Distribution over
                                             label

A C

$\theta$

                              $\theta$

$\theta$                 $\theta$            $\theta$

. . .

$pixel_{,1}$                            $pixel_{n,n}$