

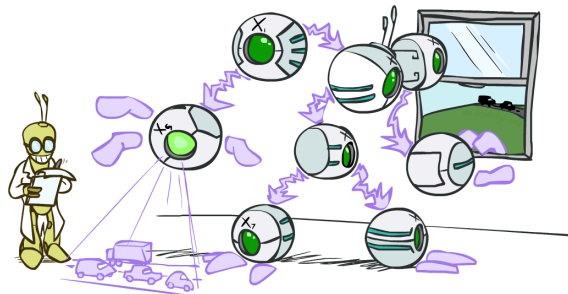
# ความรู้เบื้องต้นเกี่ยวกับปัญญาประดิษฐ์

## บรรยายครั้งที่ 5: การให้เหตุผลเชิงความน่าจะเป็น (Probabilistic Reasoning)

ผศ. ดร. อธิพล ฟองแก้ว  
[ittipon@g.sut.ac.th]

# หัวข้อในวันนี้

- โครงข่ายเบย์ (Bayesian networks)
  - ความหมาย (Semantics)
  - การสร้าง (Construction)
  - ความสัมพันธ์เชิงอิสระ (Independence relations)
- การอนุมาน (Inference)
- การเรียนรู้พารามิเตอร์ (Parameter learning)



## การแทนความรู้ที่ไม่แน่นอน

การแทนค่าการแจกแจงความน่าจะเป็นร่วม (joint probability distribution) อย่างชัดเจน จะมีขนาดเพิ่มขึ้นแบบทวีคูณ (exponentially) ตามจำนวนตัวแปร

สมมติฐานเรื่อง **ความเป็นอิสระ (Independence)** และ **ความเป็นอิสระแบบมีเงื่อนไข (conditional independence)** ช่วยลดจำนวนความน่าจะเป็นที่ต้องกำหนด เราสามารถแทนความสัมพันธ์เหล่านี้ได้อย่างชัดเจนในรูปแบบของ **โครงข่ายเบย์ (Bayesian network)**

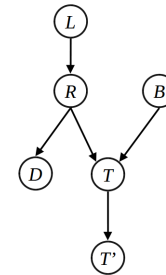
# โครงข่ายเบย์ (Bayesian networks)

โครงข่ายเบย์ คือ กราฟมีทิศทางที่ไม่มีวัฏจักร (directed acyclic graph) โดยที่

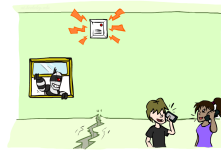
- แต่ละ โหนด (node) แทนตัวแปรสุ่ม (random variable)
  - ซึ่งอาจจะสังเกตได้ (observed) หรือสังเกตไม่ได้ (unobserved)
  - และอาจเป็นได้ทั้งแบบไม่ต่อเนื่อง (discrete) หรือต่อเนื่อง (continuous)
- แต่ละ เส้นเชื่อม (edge) จะมีทิศทางและบ่งชี้ถึงความสัมพันธ์เชิงความน่าจะเป็นโดยตรงระหว่างสองตัวแปร
- แต่ละโหนด  $X_i$  จะถูกกำกับด้วย การแจกแจงความน่าจะเป็นแบบมีเงื่อนไข (conditional probability distribution)

$$P(X_i | \text{parents}(X_i))$$

ซึ่งกำหนดการแจกแจงของ  $X_i$  เมื่อกำหนดค่าของโหนดพ่อแม่ (parents) ของมันในโครงข่าย

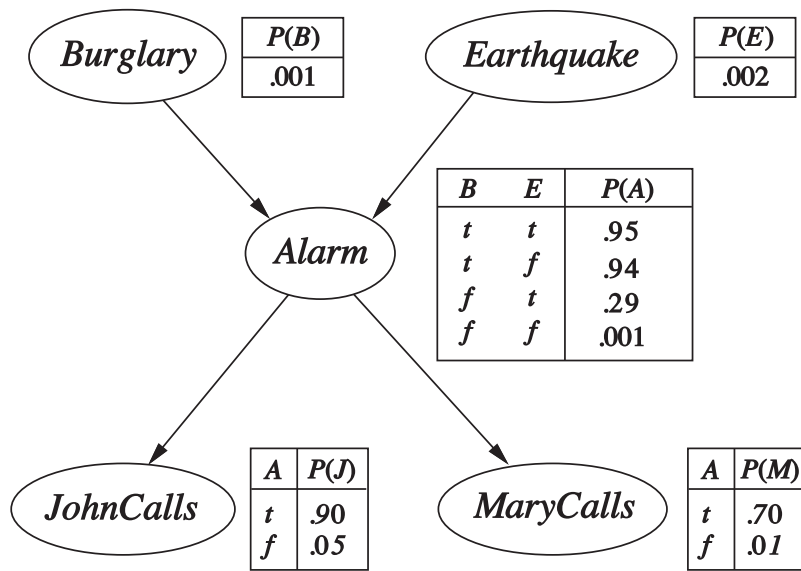


ในกรณีที่ง่ายที่สุด การแจกแจงแบบมีเงื่อนไขจะถูกแสดงในรูปแบบตารางความน่าจะเป็นแบบมีเงื่อนไข (Conditional Probability Tables - CPTs)



## ตัวอย่างที่ 1

- ตัวแปร: **Burglar** (ขโมย), **Earthquake** (แผ่นดินไหว), **Alarm** (สัญญาณเตือน), **JohnCalls** (จอห์นโทร), **MaryCalls** (แมรีโทร)
- โครงสร้างของเครือข่ายสามารถกำหนดได้จากความรู้ในสถานการณ์นั้นๆ:
  - ขโมยสามารถทำให้สัญญาณเตือนดังได้
  - แผ่นดินไหวสามารถทำให้สัญญาณเตือนดังได้
  - สัญญาณเตือนอาจทำให้แมรีโทรมา
  - สัญญาณเตือนอาจทำให้จอห์นโทรมา สถานการณ์: ฉันทูอยู่ที่ทำงาน เพื่อนบ้านชื่อจอห์นโทรมาบอกว่าสัญญาณเตือนที่บ้านฉันดัง แต่เพื่อนบ้านอีกคนชื่อแมรีไม่ได้โทรมา บางครั้งสัญญาณเตือนก็ดังจากแผ่นดินไหวเล็กน้อย คำถามคือ มีขโมยเข้าบ้านหรือไม่?



## ความหมาย (Semantics)

โครงข่ายเบย์เข้ารหัสการแจกแจงความน่าจะเป็นร่วมทั้งหมด (full joint distribution) โดยปริยาย ในรูปแบบผลคูณของการแจกแจงเฉพาะที่ (local distributions) นั่นคือ

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$

การพิสูจน์:

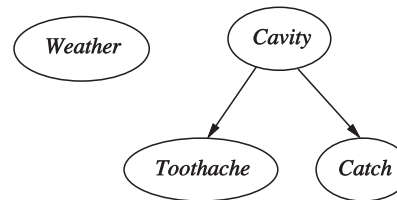
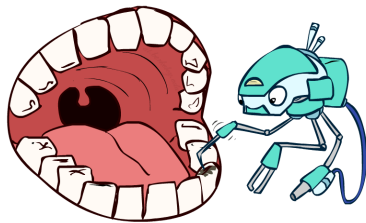
- จากกฎลูกโซ่ (chain rule),  $P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid x_1, \dots, x_{i-1})$
- หากเราตั้งสมมติฐานว่า  $X_i$  เป็นอิสระแบบมีเงื่อนไข (conditionally independent) จาก โหนดรุ่นก่อนหน้า ในลำดับที่กำหนดเมื่อให้ค่าของ โหนด  $\text{parents}(X_i) \subseteq \{X_1, \dots, X_{i-1}\}$  เราจะได้ว่า  $P(x_i \mid x_1, \dots, x_{i-1}) = P(x_i \mid \text{parents}(X_i))$
- ดังนั้น,  $P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$



**ตัวอย่างที่ 1 (ต่อ)**

$$P(j, m, a, \neg b, \neg e) = P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) = 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \approx 0.00063$$

## ตัวอย่างที่ 2



สถานการณ์ของทันตแพทย์สามารถจำลองเป็นโครงข่ายเบย์ที่มี 4 ตัวแปร ดังที่แสดงทางด้านขวา

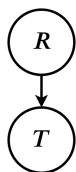
โดยโครงสร้างของเครือข่าย จะเข้ารหัสสมมติฐานความเป็นอิสระแบบมีเงื่อนไขไว้ในตัว แต่ละตัวแปรจะเป็นอิสระจากโหนดที่ไม่ใช่ลูกหลาน (non-descendants) ของมัน เมื่อกำหนดค่าของโหนดพ่อแม่:

- **Weather** (สภาพอากาศ) เป็นอิสระจากตัวแปรอื่นๆ
- **Toothache** (ปวดฟัน) และ **Catch** (ตรวจพบ) เป็นอิสระแบบมีเงื่อนไขจากกัน เมื่อกำหนดค่า **Cavity** (ฟันผุ) ทันตแพทย์กำลังตรวจฟันของคนไข้ คนไข้มีฟันผุ แต่ทันตแพทย์ยังไม่ทราบ อย่างไรก็ตาม คนไข้มีอาการปวดฟัน ซึ่งทันตแพทย์สังเกตได้



### ตัวอย่างที่ 3

เส้นเชื่อม (Edges) อาจสอดคล้องกับความสัมพันธ์เชิงสาเหตุ (causal relations)



$P(R)$

R	P
r	0.25
$\neg r$	0.75

$P(T|R)$

R	T	P
r	t	0.75
r	$\neg t$	0.25
$\neg r$	t	0.5
$\neg r$	$\neg t$	0.5

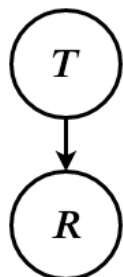
แบบจำลองเชิงสาเหตุ (Causal model): ฝนตก (R) ทำให้รถติด (T).

Courtesy: CS 168, UC Berkeley.



### ตัวอย่างที่ 3 (อีกแบบ)

... แต่เส้นเชื่อมไม่จำเป็นต้องเป็นความสัมพันธ์เชิงสาเหตุเสมอไป!



$P(T)$

T	P
t	9/16
¬t	7/16

$P(R|T)$

T	R	P
t	r	1/3
t	¬r	2/3
¬t	r	1/7
¬t	¬r	6/7

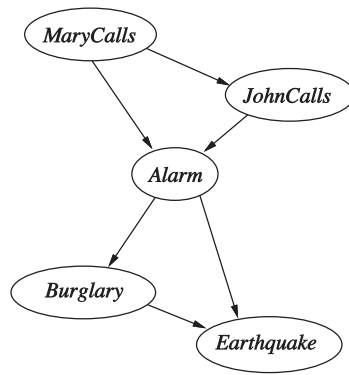
Credits: CS188, UC Berkeley.

## การสร้าง (Construction)

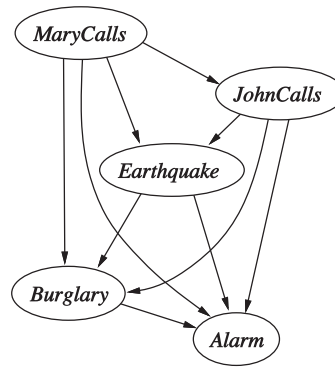
เราสามารถสร้างโครงข่ายเบย์โดยใช้ลำดับของตัวแปรแบบใดก็ได้ トラバタイยังเคารพสมมติฐานความเป็นอิสระแบบมีเงื่อนไข

### อัลกอริทึม

1. เลือก ลำดับ (ordering) ของตัวแปร  $X_1, \dots, X_n$
2. สำหรับ  $i = 1$  ถึง  $n$ :
  1. เพิ่ม  $X_i$  เข้าไปในเครือข่าย
  2. เลือกเซตของโหนดพ่อแม่ (parents) ที่เล็กที่สุดจาก  $X_1, \dots, X_{i-1}$  ที่ทำให้  $P(x_i | x_1, \dots, x_{i-1}) = P(x_i | \text{parents}(X_i))$
  3. สำหรับพ่อแม่แต่ละโหนด ให้ลากเส้นเชื่อมจากพ่อแม่ไปยัง  $X_i$
  4. สร้างตาราง CPT (Conditional Probability Table)



(a)



(b)

โครงข่ายทั้งสองนี้แทนการแจกแจงเดียวกันหรือไม่? และมีความกระชับ (compact) เท่ากันหรือไม่?

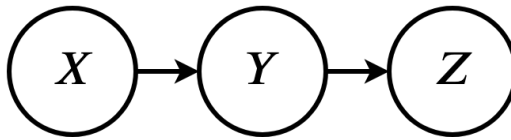
สำหรับโครงข่ายด้านซ้าย:

- $P(J|M) = P(J)$ ? ไม่
- $P(A|J, M) = P(A|J)$ ?  $P(A|J, M) = P(A)$ ? ไม่
- $P(B|A, J, M) = P(B|A)$ ? ใช่
- $P(B|A, J, M) = P(B)$ ? ไม่
- $P(E|B, A, J, M) = P(E|A)$ ? ไม่
- $P(E|B, A, J, M) = P(E|A, B)$ ? ใช่

โครงข่ายด้านขวา (โครงข่ายเดิม) มีความกระชับกว่ามาก เพราะต้องการพารามิเตอร์ใน CPT น้อยกว่า

## ความสัมพันธ์เชิงอิสระ (Independence relations)

เนื่องจากโครงสร้าง (topology) ของโครงข่ายเบย์ซาร์หัสสมมติฐานความเป็นอิสระแบบมีเงื่อนไข เราจึงสามารถเชื่อมั่นเพื่อตอบคำถามเกี่ยวกับความเป็นอิสระของตัวแปรต่างๆ เมื่อมีข้อมูลบางอย่าง (evidence)



ตัวอย่าง:  $X$  และ  $Z$  จำเป็นต้องเป็นอิสระจากกันหรือไม่?



## แนวคิดเบื้องหลัง d-separation

d-separation (ย่อมาจาก direction-dependent separation) เป็นวิธีการตรวจสอบว่ากลุ่มของตัวแปร  $X$  เป็นอิสระแบบมีเงื่อนไข (conditionally independent) จากกลุ่มของตัวแปร  $Y$  หรือไม่ เมื่อกำหนดกลุ่มของตัวแปรที่เป็นหลักฐาน (evidence)  $Z$  ให้

แนวคิดหลักคือการพิจารณา **เส้นทาง (path)** ทั้งหมดระหว่างโหนดใน  $X$  และโหนดใน  $Y$

- หากทุกเส้นทางถูก '**ปิดกั้น**' (blocked) โดยหลักฐาน  $Z$  เราจะสรุปได้ว่า  $X$  และ  $Y$  เป็นอิสระจากกันเมื่อมี  $Z$
- หากมีอย่างน้อยหนึ่งเส้นทางที่ '**เปิด**' (active) แสดงว่าอิทธิพลอาจส่งผ่านได้ และเราไม่สามารถรับประกันความเป็นอิสระได้

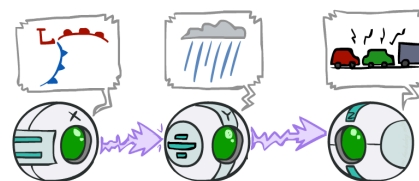
เส้นทางจะถูกพิจารณาว่าเป็น 'active' (ข้อมูลไหลผ่านได้) หรือ 'inactive' (ถูกปิดกั้น) ขึ้นอยู่กับโครงสร้างของเส้นทางและโหนดหลักฐานที่อยู่บนเส้นทางนั้น

## Cascades (แบบอนุกรม)

$X$  เป็นอิสระจาก  $Z$  หรือไม่? **ไม่**

ตัวอย่างแย้ง:

- ความกดอากาศต่ำ ( $X$ ) ทำให้เกิดฝน ( $Y$ ), ฝนทำให้รถติด ( $Z$ )
- ในทางตัวเลข:
  - $P(y|x) = 1, P(z|y) = 1$
  - $P(\neg y|\neg x) = 1, P(\neg z|\neg y) = 1$
- ข้อมูลเกี่ยวกับ  $X$  ส่งผลต่อความเชื่อของเราเกี่ยวกับ  $Z$



$X$ : ความกดอากาศต่ำ  $Y$ : ฝน  $Z$ : รถติด

$$P(x, y, z) = P(x)P(y|x)P(z|y)$$

$X$  เป็นอิสระจาก  $Z$  เมื่อกำหนด  $Y$  หรือไม่? **ใช่**

$$P(z|x,y) = \frac{P(x,y,z)}{P(x,y)} = \frac{P(x)P(y|x)P(z|y)}{P(x)P(y|x)} = P(z|y)$$

เราล่าว่าหลักฐาน (evidence)  $Y$  ตามเส้นทางแบบอนุกรมนี้  
**ปิดกั้น (blocks)** อิทธิพล



$X$ : ความกดอากาศต่ำ  $Y$ : ฝน  $Z$ : รถติด

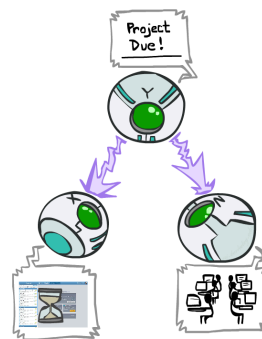
$$P(x,y,z) = P(x)P(y|x)P(z|y)$$

## Common parent (พ่อแม่ร่วม)

$X$  เป็นอิสระจาก  $Z$  หรือไม่? **ไม่**

ตัวอย่างเช่น:

- มีโปรเจกต์ที่ต้องส่ง ( $Y$ ) ทำให้ฟอร์มคนเยอะ ( $X$ ) และห้องแล็บเต็ม ( $Z$ )
- ในทางตัวเลข:
  - $P(x|y) = 1, P(\neg x|\neg y) = 1$
  - $P(z|y) = 1, P(\neg z|\neg y) = 1$
- การสังเกตว่าฟอร์มคนเยอะ ( $X$ ) ทำให้เราเชื่อมั่นมากขึ้นว่ามีโปรเจกต์ต้องส่ง ( $Y$ ) ซึ่งก็ทำให้เราเชื่อว่าห้องแล็บจะเต็ม ( $Z$ )



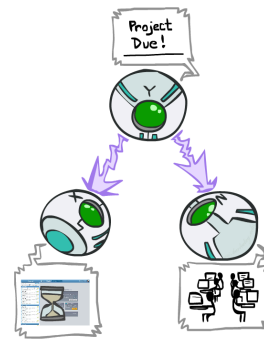
$X$ : ฟอร์มคนเยอะ  $Y$ : โปรเจกต์ต้องส่ง  $Z$ : ห้องแล็บเต็ม

$$P(x, y, z) = P(y)P(x|y)P(z|y)$$

$X$  เป็นอิสระจาก  $Z$  เมื่อกำหนด  $Y$  หรือไม่? ใช่

$$P(z|x, y) = \frac{P(x, y, z)}{P(x, y)} = \frac{P(y)P(x|y)P(z|y)}{P(y)P(x|y)} = P(z|y)$$

การสังเกตโน้ตพ่อแม่ ( $Y$ ) จะปิดกั้นอิทธิพลระหว่างโน้ตลูก ( $X, Z$ )



$X$ : ฟอรัมคนเยอะ  $Y$ : โปรเจกต์ต้องส่ง  $Z$ : ห้องแล็บเต็ม

$$P(x, y, z) = P(y)P(x|y)P(z|y)$$

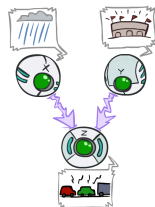
## v-structures (โครงสร้างรูปตัว V)

**X** และ **Y** เป็นอิสระจากกันหรือไม่? **ใช่**

- ฝนตก (**X**) และการแข่งขันเบสบอล (**Y**) ทำให้รถติด (**Z**) แต่การเกิดฝนกับการมีแข่งเบสบอลเป็นอิสระจากกัน
- (พิสูจน์!)

**X** และ **Y** เป็นอิสระจากกันเมื่อกำหนด **Z** หรือไม่? **ไม่!**

- การที่เราทราบว่ารถติด (**Z**) ทำให้ฝนตก (**X**) และการแข่งเบสบอล (**Y**) กลายเป็นคำอธิบายที่ "แข่งขัน" กันเอง (explaining away)  
เช่น ถ้ารู้ว่ารถติดแต่ไม่มีแข่งเบสบอล ก็จะทำให้เราเชื่อมากขึ้นว่าฝนตก
- นี่คือการที่ **ตรงกันข้าม** กับกรณีก่อนหน้านี้ การสังเกตโดนดลูกจะ **เปิด (activates)** อิทธิพลระหว่างโนดพ่อแม่ ] ]



**X**: ฝน **Y**: แข่งเบสบอล **Z**: รถติด

$$P(x, y, z) = P(x)P(y)P(z|x, y)$$

Credits: CS188, UC Berkeley.

พิสูจน์:  $P(x, y, z) = P(x)P(y)P(z|x, y)$  และ  $P(x, y, z) = P(x, y)P(z|x, y)$  ดังนั้น  $P(x, y) = P(x)P(y)$

### d-separation

สมมติว่าเรามีโครงข่ายเบย์ที่สมบูรณ์  $X_i$  และ  $X_j$  เป็นอิสระแบบมีเงื่อนไขหรือไม่ เมื่อกำหนดหลักฐาน  $Z_1 = z_1, \dots, Z_m = z_m$ ?

พิจารณาเส้นทาง (ที่ไม่กำหนดทิศทาง) ทั้งหมดจาก  $X_i$  ไปยัง  $X_j$ :

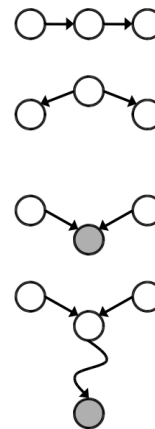
- ถ้ามีเส้นทางที่ **active** (เปิด) อย่างน้อยหนึ่งเส้นทาง จะไม่สามารถรับประกันความเป็นอิสระได้
- หากทุกเส้นทางเป็น **inactive** (ปิดกั้น) จะสามารถรับประกันความเป็นอิสระได้



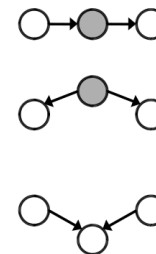
เส้นทางจะ **active** (เปิด) ถ้าแต่ละสามโหนด (triple) ตลอดเส้นทางนั้น active:

- Cascade  $A \rightarrow B \rightarrow C$  โดยที่  $B$  ไม่ถูกสังเกต (unobserved) (ไม่ว่าจะทิศทางใด)
- Common parent  $A \leftarrow B \rightarrow C$  โดยที่  $B$  ไม่ถูกสังเกต
- v-structure  $A \rightarrow B \leftarrow C$  โดยที่  $B$  หรือลูกหลาน (descendent) ของมัน ถูกสังเกต (observed)

Active Triples

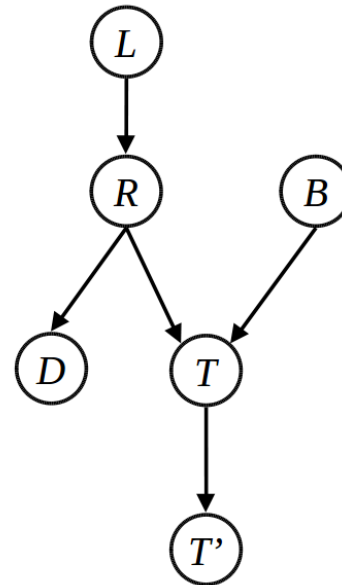


Inactive Triples



### ตัวอย่าง

- $L \perp T' | T?$
- $L \perp B?$
- $L \perp B | T?$
- $L \perp B | T'?$
- $L \perp B | T, R?$



- Yes (เส้นทาง  $L \rightarrow T \rightarrow T'$  ถูกปิดกั้นโดย  $T$ )
- Yes (เส้นทาง  $L \rightarrow T \leftarrow R \rightarrow B$  เป็น v-structure ที่  $T$  ไม่ถูกสังเกต จึงปิดกั้น)
- No (การสังเกต  $T$  ทำให้เส้นทาง  $L \rightarrow T \leftarrow R$  เปิดออก)
- No (การสังเกต  $T'$  ซึ่งเป็นลูกของ  $T$  ทำให้เส้นทาง  $L \rightarrow T \leftarrow R$  เปิดออก)
- Yes (เส้นทาง  $L \rightarrow T \leftarrow R \rightarrow B$  ถูกปิดกั้นโดย  $R$  ซึ่งเป็น common parent ที่ถูกสังเกต)

## การอนุมาน (Inference)

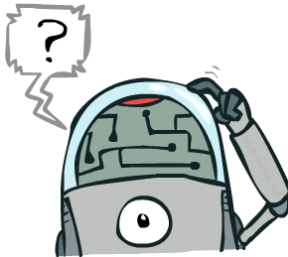
การอนุมาน (Inference) คือกระบวนการในการ **คำนวณการแจกแจงความน่าจะเป็นแบบ marginal และ/หรือ conditional** จาก **การแจกแจงความน่าจะเป็นร่วม** (joint probability distribution):

คำถามอย่างง่าย:  $P(X_i|e)$  (ความน่าจะเป็นของ  $X_i$  เมื่อมีหลักฐาน  $e$ )

คำถามแบบเชื่อม:  $P(X_i, X_j|e) = P(X_i|e)P(X_j|X_i, e)$

คำอธิบายที่เป็นไปได้ที่สุด (Most Likely Explanation):  $\arg \max_q P(q|e)$  (หาชุดค่า  $q$  ที่มีความน่าจะเป็นสูงสุดเมื่อมี  $e$ )

การตัดสินใจที่เหมาะสมที่สุด  $\arg \max_a \mathbb{E}_{p(s'|s,a)} [V(s')]$



อธิบายว่า **arg max** หมายถึงอะไร (การหาอาร์กิวเมนต์ที่ทำให้ฟังก์ชันมีค่าสูงสุด)

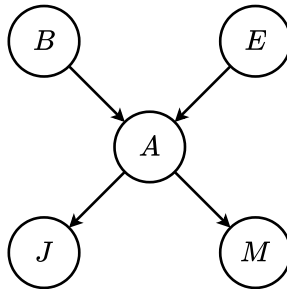
เน้นย้ำความสำคัญของการอนุมาน: การอนุมาน  $\Leftrightarrow$  การให้เหตุผล

## การอนุมานโดยการแจกแจง (Inference by enumeration)

เริ่มต้นจากการแจกแจงร่วม  $\mathbf{P}(\mathbf{Q}, \mathbf{E}_1, \dots, \mathbf{E}_k, \mathbf{H}_1, \dots, \mathbf{H}_r)$  (Q: Query, E: Evidence, H: Hidden)

1. เลือกเฉพาะรายการที่สอดคล้องกับหลักฐาน (evidence)  $\mathbf{E}_1, \dots, \mathbf{E}_k = \mathbf{e}_1, \dots, \mathbf{e}_k$
2. ทำการ Marginalize out (บวกรวมข้าม) ตัวแปรซ่อนเร้น (hidden variables) เพื่อให้ได้การแจกแจงร่วมของตัวแปรที่สนใจ (query) และตัวแปรหลักฐาน:  $\mathbf{P}(\mathbf{Q}, \mathbf{e}_1, \dots, \mathbf{e}_k) = \sum_{\mathbf{h}_1, \dots, \mathbf{h}_r} \mathbf{P}(\mathbf{Q}, \mathbf{h}_1, \dots, \mathbf{h}_r, \mathbf{e}_1, \dots, \mathbf{e}_k)$
3. ทำให้เป็นปกติ (Normalize):

$$\mathbf{Z} = \sum_{\mathbf{q}} \mathbf{P}(\mathbf{q}, \mathbf{e}_1, \dots, \mathbf{e}_k) \mathbf{P}(\mathbf{Q} | \mathbf{e}_1, \dots, \mathbf{e}_k) = \frac{1}{\mathbf{Z}} \mathbf{P}(\mathbf{Q}, \mathbf{e}_1, \dots, \mathbf{e}_k)$$



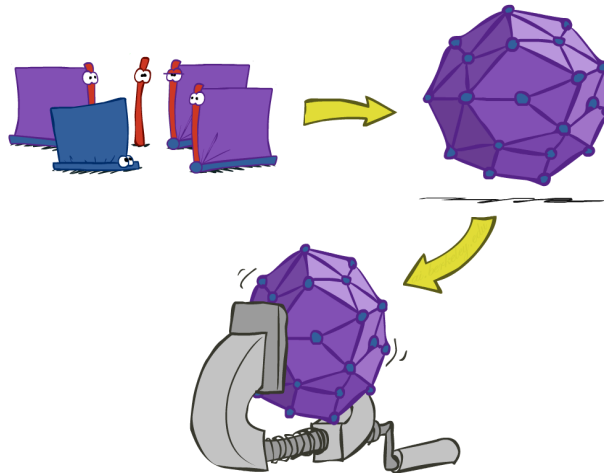
พิจารณาโครงข่ายสัญญาณเตือนและคำถาม  $\mathbf{P(B|j, m)}$  (ความน่าจะเป็นที่จะมีขโมย เมื่อจอห์นและแมรีโทรมา)

$$\mathbf{P(B|j, m)} = \frac{1}{Z} \sum_e \sum_a \mathbf{P(B, j, m, e, a)} \propto \sum_e \sum_a \mathbf{P(B, j, m, e, a)}$$

เมื่อใช้โครงข่ายเบย์ เราสามารถเขียนการแจกแจงร่วมในรูปผลคูณของค่าจาก CPT ได้:

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





การอนูมานโดยการแจกแจงนั้นช้า เพราะต้องสร้างการแจกจ่ายร่วมทั้งหมดขึ้นมาก่อนที่จะบวกรวมตัวแปรซ่อนเร้นออกไป

เราสามารถดึงตัวคูณ (factors) ที่ไม่ขึ้นกับตัวแปรในผลรวมออกมาได้ ซึ่งหมายความว่า การทำ marginalization ไม่จำเป็นต้องทำตอนท้ายสุดเสมอไป และช่วยลดการคำนวณลงได้

สำหรับโครงข่ายสัญญาณเตือน เราจะได้:

$$\mathbf{P}(B|j, m) \propto \sum_e \sum_a \mathbf{P}(B) \mathbf{P}(e) \mathbf{P}(a|B, e) \mathbf{P}(j|a) \mathbf{P}(m|a) = \mathbf{P}(B) \sum_e \mathbf{P}(e) \sum_a \mathbf{P}(a|B, e) \mathbf{P}(j|a) \mathbf{P}(m|a)$$

```

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 Bayes net with variables  $\{X\} \cup \mathbf{E} \cup \mathbf{Y}$   /*  $\mathbf{Y}$  = hidden variables */

 $Q(X) \leftarrow$  a distribution over  $X$ , initially empty
for each value  $x_i$  of  $X$  do
     $Q(x_i) \leftarrow$  ENUMERATE-ALL( $bn.VARS, \mathbf{e}_{x_i}$ )
    where  $\mathbf{e}_{x_i}$  is  $\mathbf{e}$  extended with  $X = x_i$ 
return NORMALIZE( $Q(X)$ )

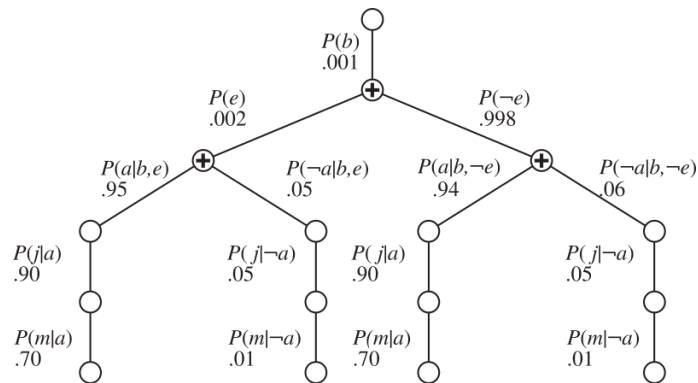
function ENUMERATE-ALL( $vars, \mathbf{e}$ ) returns a real number
if EMPTY?( $vars$ ) then return 1.0
 $Y \leftarrow$  FIRST( $vars$ )
if  $Y$  has value  $y$  in  $\mathbf{e}$ 
    then return  $P(y \mid \text{parents}(Y)) \times \text{ENUMERATE-ALL}(\text{REST}(vars), \mathbf{e})$ 
    else return  $\sum_y P(y \mid \text{parents}(Y)) \times \text{ENUMERATE-ALL}(\text{REST}(vars), \mathbf{e}_y)$ 
    where  $\mathbf{e}_y$  is  $\mathbf{e}$  extended with  $Y = y$ 

```

มีความซับซ้อนเช่นเดียวกับ DFS:  $O(n)$  ในด้านพื้นที่ และ  $O(d^n)$  ในด้านเวลา

- $n$  คือจำนวนตัวแปร
- $d$  คือขนาดของโดเมนของตัวแปร

แผนภูมิต้นไม้การประเมินสำหรับ  $P(b|j, m)$



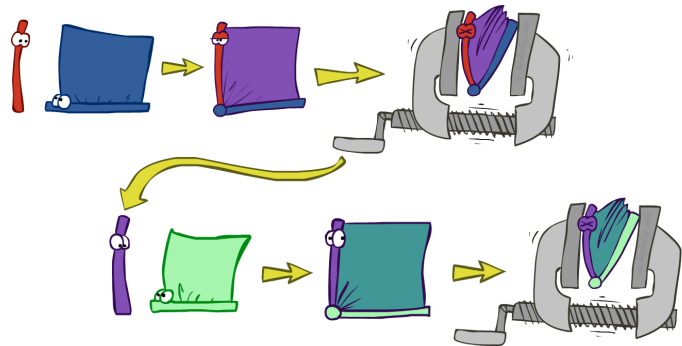
แม้จะมีการแยกตัวคุณออกมา การอนุมาณโดยการแจกแจงก็ยังคง **ไม่มีประสิทธิภาพ** เพราะมีการคำนวณซ้ำซ้อน!

- เช่น  $P(j|a)P(m|a)$  ถูกคำนวณสองครั้ง ครั้งหนึ่งสำหรับ  $e$  และอีกครั้งสำหรับ  $\neg e$
- เราสามารถหลีกเลี่ยงปัญหานี้ได้โดยการเก็บ **ผลลัพธ์ระหว่างกลาง** ไม่มีประสิทธิภาพเพราะผลคูณถูกประเมินจากซ้ายไปขวาใน **ลักษณะของ DFS**

## การอนุมานโดยการกำจัดตัวแปร (Variable Elimination)

อัลกอริทึม **Variable Elimination** จะทำการบวกรวมจากขวาไปซ้าย และเก็บตัวคุณระหว่งกลาง (intermediate factors) เพื่อหลีกเลี่ยงการคำนวณซ้ำ อัลกอริทึมจะทำงานสลับกันระหว่าง:

- การรวมตารางย่อย (Joining sub-tables)
- การกำจัดตัวแปรซ่อนเร้น (Eliminating hidden variables)



## Variable Elimination

คำถาม:  $P(Q|e_1, \dots, e_k)$

1. เริ่มต้นด้วยตัวคูณเริ่มต้น (initial factors) (คือ CPTs เฉพาะที่ ซึ่งถูกแทนค่าด้วยหลักฐานแล้ว)
2. トラバídaที่ยังมีตัวแปรซ่อนเร้น:
  1. เลือกตัวแปรซ่อนเร้น  $H$
  2. รวม (Join) ตัวคูณทั้งหมดที่เกี่ยวข้องกับ  $H$
  3. กำจัด (Eliminate)  $H$
3. รวมตัวคูณที่เหลือทั้งหมด
4. ทำให้เป็นปกติ (Normalize)

### ตัวคูณ (Factors)

- ตัวคูณ  $f_i$  แต่ละตัวคืออาร์เรย์หลายมิติที่ถูกทำดัชนีด้วยค่าของตัวแปรที่เป็นอาร์กิวเมนต์ของมัน เช่น:

$$f_4 = f_4(A) = (P(j|a) \ P(j|\neg a)) = (0.90 \ 0.05) \quad f_4(a) = 0.90 \quad f_4(\neg a) = 0.05$$

- ตัวคูณจะถูกกำหนดค่าเริ่มต้นด้วย CPTs ที่กำกับโหนดของโครงข่ายเบย์ โดยปรับตามหลักฐาน (evidence) ที่มี

## การรวม (Join)

ผลคูณแบบจุดต่อจุด (pointwise product)  $\times$ , หรือ  $\text{*join}$ , ของสองตัวคูณ  $f_1$  และ  $f_2$  จะได้ตัวคูณใหม่  $f_3$

- เหมือนกับการทำ **database join**!
- ตัวแปรของ  $f_3$  คือ **ยูเนียน** ของตัวแปรใน  $f_1$  และ  $f_2$
- สมาชิกของ  $f_3$  ได้มาจากการคูณสมาชิกที่สอดคล้องกันใน  $f_1$  และ  $f_2$

$A$	$B$	$f_1(A, B)$	$B$	$C$	$f_2(B, C)$	$A$	$B$	$C$	$f_3(A, B, C)$
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:  $f_1(A, B) \times f_2(B, C) = f_3(A, B, C)$ .

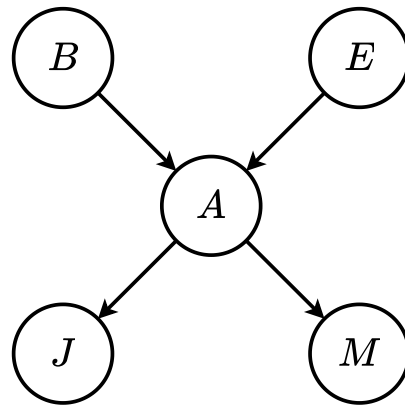


### การกำจัด (Elimination)

การบวกรวมข้าม (Summing out), หรือ การกำจัด (eliminating), ตัวแปรออกจากตัวคูณ ทำได้โดยการบวกรวมอาร์เรย์ย่อยที่เกิดจากการกำหนดค่าของตัวแปรนั้นให้เป็นแต่ละค่าที่เป็นไปได้

ตัวอย่างเช่น เพื่อกำจัด  $A$  ออกจาก  $f_3(A, B, C)$ , เราจะเขียนว่า:

$$f(B, C) = \sum_a f_3(a, B, C) = f_3(a, B, C) + f_3(\neg a, B, C) = (0.06 \quad 0.24 \quad 0.42 \quad 0.28) + (0.18 \quad 0.72 \quad 0.06 \quad 0.04) = (0.24 \quad 0.96 \quad 0.48 \quad 0.32)$$



ลองใช้วิธีการ variable elimination สำหรับคำถาม  $P(B|j, m)$

## ความเกี่ยวข้อง (Relevance)

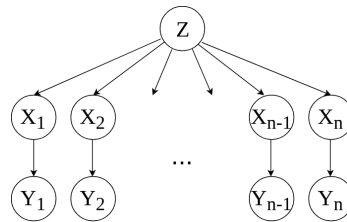
พิจารณาคำถาม  $P(J|b)$ :

$$P(J|b) \propto P(b) \sum_e P(e) \sum_a P(a|b, e) P(J|a) \sum_m P(m|a)$$

- $\sum_m P(m|a) = 1$  (เพราะผลรวมความน่าจะเป็นของทุกค่าที่เป็นไปได้ของ  $M$  ต้องเท่ากับ 1) ดังนั้น  $M$  จึง **ไม่เกี่ยวข้อง** (irrelevant) กับคำถามนี้
- กล่าวอีกนัยหนึ่ง  $P(J|b)$  จะไม่เปลี่ยนแปลงถ้าเรานำ  $M$  ออกจากเครือข่าย

ทฤษฎีบท. ตัวแปร  $H$  จะไม่เกี่ยวข้อง กับ  $P(Q|e)$  เว้นแต่ว่า  $H$  จะอยู่ในเซตของบรรพบุรุษ (ancestors) ของ  $Q \cup E$

## ความซับซ้อน (Complexity)



พิจารณาคำถาม  $P(X_n | y_1, \dots, y_n)$

ลองพิจารณาลำดับการกำจัดสองแบบ:

- $SZ, X_1, \dots, X_{\{n-1\}}S$
- $SX_1, \dots, X_{\{n-1\}}, ZS$

ขนาดของตัวคูณที่ใหญ่ที่สุดที่สร้างขึ้นสำหรับแต่ละลำดับคือเท่าใด?

- คำตอบ:  $2^{n+1}$  เทียบกับ  $2^2$  (สมมติว่าเป็นตัวแปรแบบบูลีน)

ความซับซ้อนทางเวลาและพื้นที่ของ variable elimination ถูกกำหนดโดยขนาดของตัวคูณที่ใหญ่ที่สุด

- ลำดับ การกำจัดตัวแปรมีส่วนอย่างมากต่อขนาดของตัวคูณที่ใหญ่ที่สุด
- การหาลำดับที่ดีที่สุดเป็นปัญหา **NP-hard** ซึ่งหมายความว่าไม่มีอัลกอริทึมที่หาคำตอบได้ในเวลาพหุนาม (polynomial-time)

## การอนุมานเชิงประมาณ (Approximate inference)

การอนุมานแบบแม่นยำตรง (Exact inference) นั้น **ไม่สามารถทำได้ในทางปฏิบัติ (intractable)** สำหรับแบบจำลองที่น่าจะเป็นส่วนใหญ่ที่น่าสนใจในโลกแห่งความเป็นจริง (เช่น แบบจำลองที่มีตัวแปรจำนวนมาก, มีทั้งตัวแปรต่อเนื่องและไม่ต่อเนื่อง, มีวงจรแบบไม่มีทิศทาง ฯลฯ)

เราจึงต้องหันไปใช้อัลกอริทึมการอนุมาน **เชิงประมาณ (approximate)**:

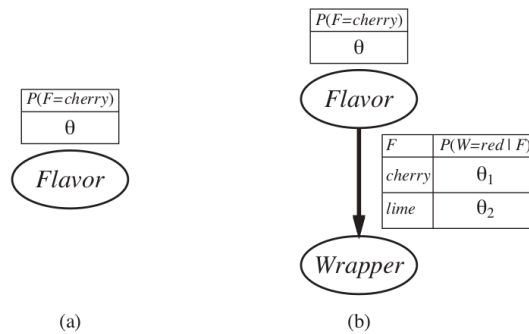
- **วิธีสุ่มตัวอย่าง (Sampling methods)**: สร้างคำตอบโดยการสุ่มตัวอย่างซ้ำๆ จากการแจกแจงที่สนใจ
- **วิธีเชิงการแปรผัน (Variational methods)**: กำหนดปัญหาการอนุมานให้เป็นปัญหาการหาค่าที่เหมาะสมที่สุด (optimization problem)
- **วิธี Belief propagation**: กำหนดการอนุมานเป็นอัลกอริทึมการส่งผ่านข้อความ (message-passing)
- **วิธี Machine learning**: เรียนรู้การประมาณค่าของการแจกแจงเป้าหมายจากข้อมูลตัวอย่าง

## การเรียนรู้พารามิเตอร์ (Parameter learning)

เมื่อสร้างแบบจำลองสำหรับโดเมนใดๆ เราสามารถเลือกแบบจำลองความน่าจะเป็นที่ระบุเป็นโครงข่ายเบย์ได้ อย่างไรก็ตาม การกำหนดค่าความน่าจะเป็นแต่ละค่าด้วยตนเองมักเป็นเรื่องยาก

ทางออกคือการใช้กลุ่มของแบบจำลองที่มีพารามิเตอร์ (parameterized)  $P(X|\theta)$  (บางครั้งเขียนเป็น  $P_\theta(X)$ ) และ ประเมินค่า (estimate) พารามิเตอร์  $\theta$  จากข้อมูล





**Figure 20.2** (a) Bayesian network model for the case of candies with an unknown proportion of cherries and limes. (b) Model for the case where the wrapper color depends (probabilistically) on the candy flavor.

## การประมาณค่าด้วยความควรจะเป็นสูงสุด (Maximum Likelihood Estimation)

สมมติว่าเรามีชุดข้อมูลสังเกต  $N$  ชุดที่เป็นอิสระและมีการแจกแจงเหมือนกัน (i.i.d.)  $\mathbf{d} = \mathbf{x}_1, \dots, \mathbf{x}_N$

ความควรจะเป็น (likelihood) ของพารามิเตอร์  $\theta$  คือความน่าจะเป็นของข้อมูลเมื่อกำหนดพารามิเตอร์

$$P(\mathbf{d}|\theta) = \prod_{j=1}^N P(\mathbf{x}_j|\theta)$$

ค่าประมาณความควรจะเป็นสูงสุด (Maximum Likelihood Estimate - MLE)

$\theta$  ของพารามิเตอร์ คือค่าของ  $\theta$  ที่ทำให้ค่า likelihood สูงที่สุด

$$\theta^* = \arg \max_{\theta} P(\mathbf{d}|\theta)$$

ในทางปฏิบัติ:

1. เขียนฟังก์ชัน log-likelihood  $L(\theta) = \log P(\mathbf{d}|\theta)$  ของพารามิเตอร์  $\theta$
2. หาอนุพันธ์  $\frac{\partial L}{\partial \theta}$  ของ log-likelihood เทียบกับพารามิเตอร์  $\theta$
3. หาค่าพารามิเตอร์  $\theta^*$  ที่ทำให้อนุพันธ์เป็นศูนย์ (และตรวจสอบว่า Hessian เป็น negative definite เพื่อให้แน่ใจว่าเป็นค่าสูงสุด)

### กรณี (a)

สัดส่วน  $\theta$  ของลูกอมรสเชอร์รี่เป็นเท่าใด?

สมมติว่าเราแกะลูกอม  $N$  เม็ด ได้รสเชอร์รี่  $c$  เม็ด และรสมะนาว  $l = N - c$  เม็ด ข้อมูลเหล่านี้เป็น i.i.d. ดังนั้น

$$P(\mathbf{d}|\theta) = \prod_{j=1}^N P(x_j|\theta) = \theta^c (1 - \theta)^l$$

หาค่าสูงสุดของฟังก์ชันนี้เทียบกับ  $\theta$  ซึ่งจะง่ายกว่าถ้าทำกับ log-likelihood:

$$L(\mathbf{d}|\theta) = \log P(\mathbf{d}|\theta) = c \log \theta + l \log(1 - \theta) \quad \frac{\partial L(\mathbf{d}|\theta)}{\partial \theta} = \frac{c}{\theta} - \frac{l}{1 - \theta} = 0$$

ดังนั้น  $\theta = \frac{c}{N}$

### กรณี (b)

เปลือกลูกอมสีแดงและสีเขียวขึ้นอยู่กับรสชาติอย่างน่าจะเป็น เช่น likelihood ของลูกอมรสเชอร์รี่ในเปลือกสีเขียวคือ

$$P(\text{cherry, green}|\theta, \theta_1, \theta_2) = P(\text{cherry}|\theta, \theta_1, \theta_2)P(\text{green}|\text{cherry}, \theta, \theta_1, \theta_2) = \theta(1 - \theta_1)$$

likelihood ของพารามิเตอร์ เมื่อมีลูกอม  $N$  เม็ด, เชอร์รี่เปลือกแดง  $r_c$  เม็ด, เชอร์รี่เปลือกเขียว  $g_c$  เม็ด, ฯลฯ คือ

$$P(\mathbf{d}|\theta, \theta_1, \theta_2) = \theta^c (1 - \theta)^l \theta_1^{r_c} (1 - \theta_1)^{g_c} \theta_2^{r_1} (1 - \theta_2)^{g_1} \quad L = c \log \theta + l \log(1 - \theta) + r_c \log \theta_1 + g_c \log(1 - \theta_1) + r_1 \log \theta_2 + g_1 \log(1 - \theta_2)$$

อนุพันธ์ของ  $L$  จะให้ผลลัพธ์:

$$\frac{\partial L}{\partial \theta} = \frac{c}{\theta} - \frac{1}{1-\theta} = 0 \Rightarrow \theta = \frac{c}{c+1} \quad \frac{\partial L}{\partial \theta_1} = \frac{r_c}{\theta_1} - \frac{g_c}{1-\theta_1} = 0 \Rightarrow \theta_1 = \frac{r_c}{r_c + g_c} \quad \frac{\partial L}{\partial \theta_2} = \frac{r_l}{\theta_2} - \frac{g_l}{1-\theta_2} = 0 \Rightarrow \theta_2$$

ในกรณี (a) ถ้าเราแกะลูกอม 1 เม็ดและได้รสเซอร์รี่ ค่า MLE คืออะไร? เรามั่นใจในค่าประมาณนี้แค่ไหน?

- เมื่อมีข้อมูลน้อย การประมาณค่าด้วย MLE อาจนำไปสู่การ **overfitting** (การปรับตัวเข้ากับข้อมูลมากเกินไป)
- MLE ไม่ได้ให้ข้อมูลเกี่ยวกับความไม่แน่นอน (uncertainty) ของพารามิเตอร์

## การเรียนรู้พารามิเตอร์แบบเบย์ (Bayesian parameter learning)

เราสามารถมองปัญหาการเรียนรู้พารามิเตอร์เป็นปัญหา การอนุมานแบบเบย์ (Bayesian inference) ได้:

- ทำให้พารามิเตอร์  $\theta$  เป็นตัวแปรสุ่ม (random variables) และมองว่าเป็นตัวแปรซ่อนเร้น
- กำหนดการแจกแจง **ก่อนหน้า (prior)**  $P(\theta)$  ให้กับพารามิเตอร์
- จากนั้น เมื่อมีข้อมูลเข้ามา เราจะอัปเดตความเชื่อของเราเกี่ยวกับพารามิเตอร์เพื่อหาการแจกแจง **ภายหลัง (posterior)**  $P(\theta | \mathbf{d})$

ในมุมมองนี้ พารามิเตอร์ไม่ใช่ค่าคงที่ที่เราพยายามจะหาค่าที่ดีที่สุดค่าเดียวอีกต่อไป แต่เป็นตัวแปรสุ่มที่มีการแจกแจงความน่าจะเป็นของตัวเอง ซึ่งจะเปลี่ยนแปลงไปเมื่อเราได้รับข้อมูลมากขึ้น



## กรณี (a)

สัดส่วน  $\theta$  ของลูกอมรสเชอร์รี่เป็นเท่าใด?

เราสมมติ prior แบบ Beta:

$$P(\theta) = \text{Beta}(\theta|a, b) = \frac{1}{Z} \theta^{a-1} (1 - \theta)^{b-1}$$

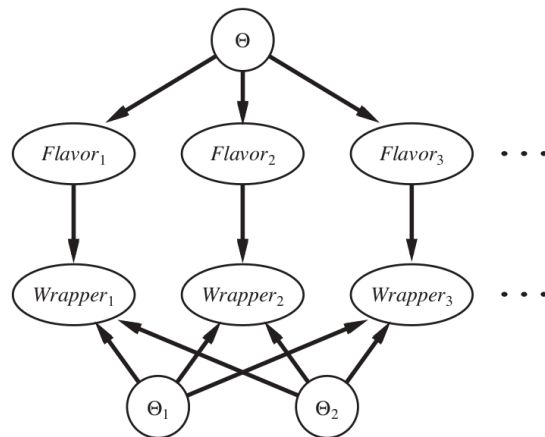
โดยที่  $Z$  คือค่าคงที่เพื่อให้เป็นปกติ (normalization constant)

จากนั้น เมื่อสังเกตเห็นลูกอมรสเชอร์รี่ จะได้ posterior:

$$P(\theta|\text{cherry}) \propto P(\text{cherry}|\theta)P(\theta) = \theta \cdot \text{Beta}(\theta|a, b) \propto \theta \cdot \theta^{a-1} (1 - \theta)^{b-1} = \theta^a (1 - \theta)^{b-1} = \text{Beta}(\theta|a + 1, b)$$

การสังเกตข้อมูลใหม่ จะเป็นการอัปเดตพารามิเตอร์ของ Beta distribution

กรณี (b)



**Figure 20.6** A Bayesian network that corresponds to a Bayesian learning process. Posterior distributions for the parameter variables  $\Theta$ ,  $\Theta_1$ , and  $\Theta_2$  can be inferred from their prior distributions and the evidence in the  $Flavor_i$  and  $Wrapper_i$  variables.

### การประมาณค่า Posterior สูงสุด (Maximum a posteriori - MAP)

เมื่อไม่สามารถคำนวณ posterior แบบ (analytically) ได้ เราสามารถใช้การประมาณค่า **Maximum a posteriori (MAP)** ซึ่งประกอบด้วยการประมาณค่า posterior ด้วยจุดประมาณ  $\theta^*$  ที่ทำให้การแจกแจง posterior มีค่าสูงสุด นั่นคือ

$$\begin{aligned} \theta^* = \arg \max_{\theta} P(\theta | \mathbf{d}) = \arg \max_{\theta} P(\mathbf{d} | \theta) P(\theta) \end{aligned}$$

MAP คล้ายกับ MLE แต่เพิ่มพจน์ของ prior  $P(\theta)$  เข้ามา ซึ่งทำหน้าที่เหมือน "regularizer" ป้องกันไม่ให้ค่าพารามิเตอร์สุดโต่งเกินไปเมื่อมีข้อมูลน้อย

## สรุป

- **โครงข่ายเบย์ (Bayesian Network)** ระบุการแจกแจงร่วมทั้งหมด และมักจะมีขนาดเล็กกว่าการแจกแจงร่วมที่แจกแจงอย่างชัดเจนแบบทวิคูณ
- **โครงสร้าง (Topology)** ของโครงข่ายเบย์เข้ารหัสสมมติฐานความเป็นอิสระแบบมีเงื่อนไขระหว่างตัวแปรสุ่ม
- **การอนุมาน (Inference)** คือปัญหาของการคำนวณการแจกแจงความน่าจะเป็นแบบ marginal และ/หรือ conditional
  - การอนุมานแบบแม่นยำตรง (Exact inference) สามารถทำได้สำหรับโครงข่ายเบย์อย่างง่าย แต่ไม่สามารถทำได้ในทางปฏิบัติสำหรับแบบจำลองส่วนใหญ่
  - อัลกอริทึมการอนุมานเชิงประมาณ (Approximate inference) ถูกนำมาใช้ในทางปฏิบัติ
- **พารามิเตอร์** ของโครงข่ายเบย์สามารถเรียนรู้ได้จากข้อมูลโดยใช้ Maximum Likelihood Estimation (MLE) หรือ Bayesian inference

จบ