

2/4/14

1

# BN PARAMETER LEARNING

## ① Setup:

→ Random Variables:  $\vec{X} = \{x_1, \dots, x_n\}$

→ "True" Distribution:  $\sim P^*(x_1, \dots, x_n)$

→ Fully Observed Dataset of samples from  $P^*$

$$\begin{bmatrix}
 x_1^{(1)} & \dots & x_n^{(1)} \\
 x_1^{(2)} & \dots & x_n^{(2)} \\
 \vdots & & \vdots \\
 x_1^{(M)} & \dots & x_n^{(M)}
 \end{bmatrix}_{M \times n}$$

→ Known BN structure is assume  $P^*$  factorizes to  $G$ .

$$P^*(x_1, \dots, x_n) = \prod_i P(x_i | \underbrace{Pa_{x_i}}_{\text{parents in } G})$$

Goal: estimate CPTs

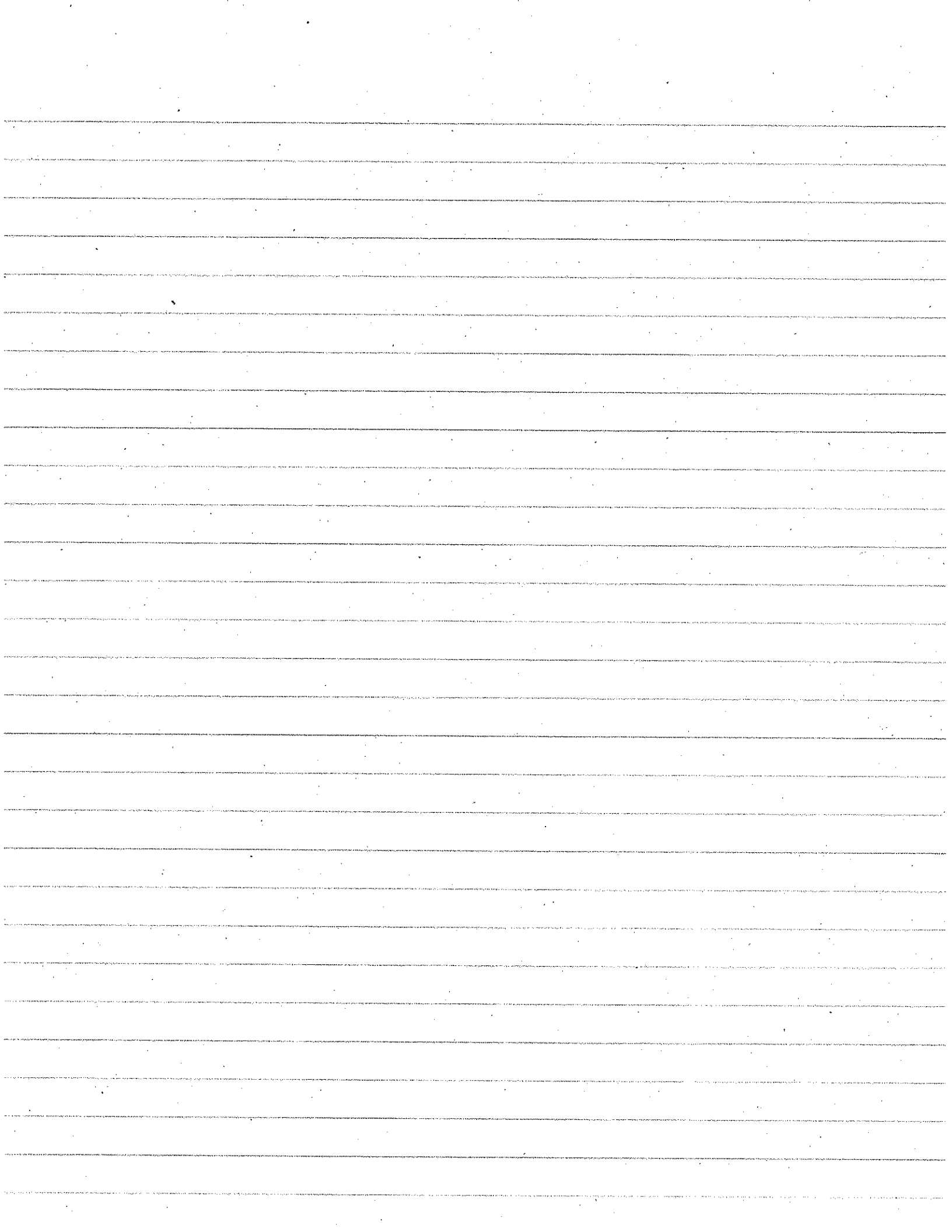
Tools: → Frequentist: Max-Likelihood

→ Bayesian: Maximum a Posteriori  
or full posterior estimation.

## ② Simplest BN: 1 RV: $(X)$

CPT  $P(X) = \begin{bmatrix} 1-\theta \\ \theta \end{bmatrix} \begin{matrix} x=0 / T \\ x=1 / F \end{matrix}$

Notes on next page from ECE 4984/5984



## ② Max Likelihood Estimation

(Sample Space)

$$\Omega = \{ \text{Nadal Loses (L)}, \text{Nadal Wins (W)} \}$$

Random Variable

$$x \in \{0, 1\}$$

L, W

T, H

$$X \sim \text{Bernoulli}(\theta)$$

$$P(X=1) = \theta$$

$$P(X=0) = 1 - \theta$$

Given:

$$D = \{1, 0, 0, 1, 1\} \quad \text{estimate } \theta \in [0, 1]$$

Good  $\theta \equiv$  makes it likely for us to have observed  $D$

[e.g. if  $\theta=0$ , we would never see 1]

Idea: let's maximize the probability of  $D$  under  $\theta$

$$\hat{\theta}_{MLE} = \underset{\theta \in [0, 1]}{\text{argmax}} \underbrace{P(D|\theta)}$$

called likelihood function

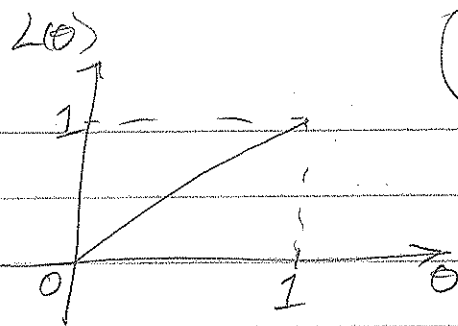
$$L(\theta; D) = P(D|\theta)$$

IMP: Likelihood is a function of  $\theta$   
( $D$  is fixed to what we observed)

$$L(\theta; D) = P(D|\theta) = \prod_{i=1}^n P(x_i|\theta) \quad [\text{Why? Hint: IID}]$$

E.g.  $D = \{1\}$

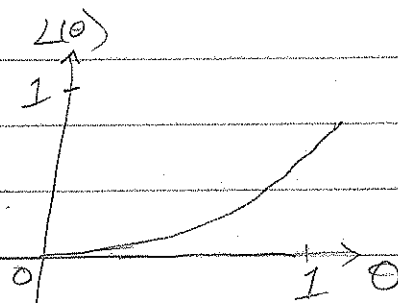
$L(\theta) = \theta$



(2)

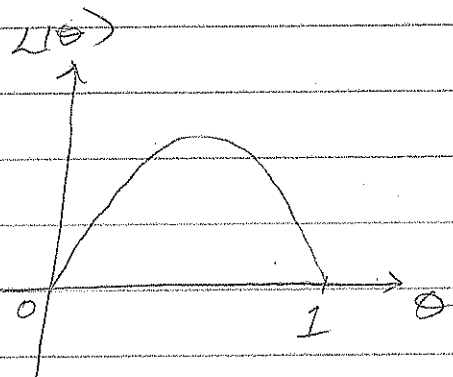
$D = \{1, 1\}$

$L(\theta) = \theta \cdot \theta = \theta^2$



$D = \{1, 0\}$

$L(\theta) = \theta \cdot (1 - \theta)$



In general  $L(\theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$

$\alpha_H = \# \text{ Heads / Wins}$

$\alpha_T = \# \text{ Tails / Losses}$

$\hat{\theta}_{MLE} = \underset{\theta \in [0,1]}{\text{argmax}} L(\theta)$

$= \underset{\theta}{\text{argmax}} \log L(\theta)$  [why?  $\because$  log is a monotone function so preserves argmax]

$\leftarrow$  called log-likelihood

How do we find argmax of  $L(\theta)$ ?

Take 1st derivative; set to zero

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} \left[ \alpha_H \log \theta + \alpha_T \log(1 - \theta) \right]$$

$$= \frac{\alpha_H}{\theta} + \frac{\alpha_T}{1 - \theta} (-1)$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = 0 \Rightarrow \frac{\alpha_H - \alpha_H \theta - \alpha_T \theta}{\theta(1-\theta)} = 0$$

Ignoring boundary conditions  $\Rightarrow$

$$\theta_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

③ Sufficient statistic

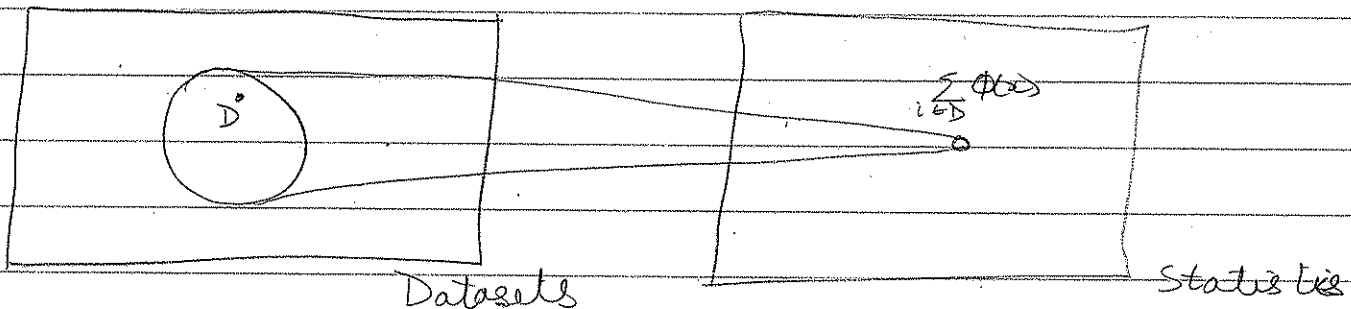
$$D_1 = \{1, 1, 1, 0, 0, 0\}$$

$$\alpha_H = \alpha_T = 3$$

$$D_2 = \{1, 0, 1, 0, 1, 0\}$$

$$\alpha_H = \alpha_T = 3$$

Two datasets but look the same to Likelihood.



④ MLE is OPT if model class is correct.

$$\text{Consider } \frac{1}{N} \mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^n \log P(X_i | \theta)$$

$$= \frac{1}{N} \left[ \#(X=1) P(X=1|\theta) + \#(X=2) P(X=2|\theta) + \dots \right]$$

(counting argument)

(3)

As data becomes infinite

$$\lim_{N \rightarrow \infty} \frac{\#(X=1)}{N} = P(X=1 | \theta^*)$$

← true parameter (unknown)

Shorthand  $\left\{ \begin{array}{l} P^*(x) = P(X | \theta^*) \\ P_\theta(x) = P(X | \theta) \end{array} \right\}$

$$\text{Now } \frac{1}{N} \text{LL}(\theta) \underset{\text{as } N \rightarrow \infty}{=} \sum_{x=1}^k P^*(x) \log P_\theta(x)$$

$$= \sum_{x=1}^k P^*(x) \log \left[ P_\theta(x) \cdot \frac{P^*(x)}{P^*(x)} \right]$$

$$= \sum_{x=1}^k P^*(x) \log P^*(x) - \sum_{x=1}^k P^*(x) \log \frac{P^*(x)}{P_\theta(x)}$$

$$\frac{1}{N} \text{LL}(\theta) = -H(P^*) - \text{KL}(P^* || P_\theta)$$

$$\text{So } \max_{\theta} \frac{1}{N} \text{LL}(\theta) \equiv \min_{\theta} \text{KL}(P^* || P_\theta)$$

[∵  $H(P^*)$  is a constant]

POWERFUL RESULT: → We did not specify  $P(X|\theta)$   
 → Any distribution

Concave: → Inf data

→ We must know the family  $P_\theta$ , which we usually don't (e.g. is life Gaussian?)

### ⑤ MAP + Bayesian Estimation

Let's think of  $\theta$  as a random quantity & apply

$$\text{Bayes Rule: } P(\theta | D) = \frac{P(D | \theta) P(\theta)}{P(D)}$$

$P(\theta)$ : What do we believe about  $\theta$  without any data?  
: Prior Belief

$\theta \in [0, 1]$  So we need a distribution over parameters of our distribution

Beta distribution:

$$P(\theta | \beta_H, \beta_T) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{\text{constant}}$$

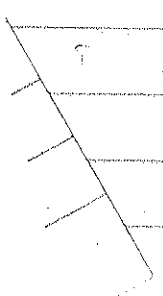
↑ ↑

hyper-parameters: parameters of the distribution over parameter ( $\theta$ )

Important Facts:

$$\text{constant} = \int_0^1 \theta^{\beta_H-1} (1-\theta)^{\beta_T-1} d\theta$$

$$\text{mode of distribution} = \frac{\beta_H - 1}{\beta_H + \beta_T - 2}$$



(4)

# Maximum A Posteriori (MAP) Estimation

$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} P(\theta | D)$$

$$= \operatorname{argmax}_{\theta} \frac{P(D | \theta) P(\theta)}{P(D)}$$

← constant w.r.t  $\theta$

$$= \operatorname{argmax}_{\theta} P(D | \theta) P(\theta)$$

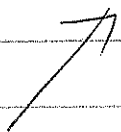
$$= \operatorname{argmax}_{\theta} \theta^{\alpha_H} (1-\theta)^{\alpha_T} \times \theta^{\beta_H-1} (1-\theta)^{\beta_T-1}$$

$$= \operatorname{argmax}_{\theta} \theta^{\alpha_H + \beta_H - 1} (1-\theta)^{\alpha_T + \beta_T - 1}$$

Beta( $\alpha_H + \beta_H$ ,  $\alpha_T + \beta_T$ )

= mode of Beta

$$= \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$



Very Nice, so  $\beta_H, \beta_T$  act as pseudo-flips/count  
 → phantom experiments not contained in our dataset



Special Cases:  $\beta_H = \beta_T = 1$

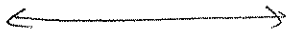
$$\hat{\theta}_{MAP} = \frac{\alpha_H + 1 - 1}{\alpha_H + 1 + \beta_T + 1 - 2} = \hat{\theta}_{MLE}$$

When  $\beta_H = \beta_T = 1$

$$\text{Beta} \equiv \theta^{\alpha} (1-\theta)^{\beta}$$

uniform distribution

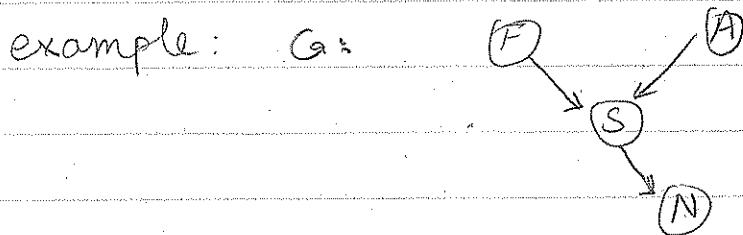
No prior  $\Rightarrow$  posterior = prior



When  $N \rightarrow \infty$  effect of  $\beta_H, \beta_T$  is forgotten.

←————→  
Back to ECE 6504 Notes.

③ All ideas extend to general BN CPT estimation



Vars:  $F, A, S, N$  (all binary)

CPTs:  $P(F) = \begin{bmatrix} 1-\theta_F \\ \theta_F \end{bmatrix}$      $P(A) = \begin{bmatrix} 1-\theta_A \\ \theta_A \end{bmatrix}$

$P(S|F, A) =$

	$F, A$	00	01	10	11
S	0	$1-c$	$1-c$	$1-c$	$1-c$
1		$\theta_{S100}$	$\theta_{S101}$	$\theta_{S110}$	$\theta_{S111}$

$P(N|S) =$

	NS	0	1
N	0	$1-l$	$1-l$
1		$\theta_{N10}$	$\theta_{N11}$

Parameters to be estimated:  $\vec{\theta} = \{ \theta_F, \theta_A, \theta_{S100}, \theta_{S101}, \theta_{S110}, \theta_{S111}, \theta_{N10}, \theta_{N11} \}$

MLE:  $\vec{\theta}_{MLE} = \underset{\vec{\theta}}{\operatorname{argmax}} \log P(D | \vec{\theta}, G)$

$= \underset{\vec{\theta}}{\operatorname{argmax}} \sum_{j=1}^M \log P(F=f^{(j)}, A=a^{(j)}, S=s^{(j)}, N=n^{(j)} | \vec{\theta}, G)$

$= \underset{\vec{\theta}}{\operatorname{argmax}} \sum_{j=1}^M \left[ \log P(F=f^{(j)} | \vec{\theta}, G) + \log P(A=a^{(j)} | \vec{\theta}, G) \right.$   
 $\left. + \log P(S=s^{(j)} | F=f^{(j)}, A=a^{(j)}, \vec{\theta}, G) + \log P(N=n^{(j)} | S=s^{(j)}, \vec{\theta}, G) \right]$

independent subproblems  
each term depends only on its own  $\theta$

same argument as for a single variable

In general

$$\theta_{x_i=a | P_{x_i=B}} = \frac{\text{Count}(X_i=a, P_{x_i=B})}{\text{Count}(P_{x_i=B})}$$

[For multi-label variables  $x_i \in \{1, \dots, k\}$ , need just a bit more math (Lagrange Multipliers)  $\sum_a \theta_{x_i=a | P_{x_i=B}} = 1$  constrained optimization; But same result as binary vars.]

For MAP estimation

assume  $P(\theta_{x_i | P_{x_i=B}}) \sim \text{Dir}(\vec{\alpha}_{x_i | P_{x_i=B}})$

Each CPT column has it's own prior given by hyperparameters