

Learning Models with Shared Parameters

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Overview
 1. Global Parameter Sharing
 2. Local Parameter Sharing
 3. Bayesian Inference with Shared Parameters
 4. Hierarchical Priors

Overview

- Parameters of different CPDs in a BN are not necessarily independent
 - Same parameters occur in multiple places in a BN
- Methods to estimate shared parameters
 - MLE methods that exploit local and global structure
 1. Global : same CPD used across multiple variables
 - Occurs naturally in template-based models
 2. Local: finer grained
 - Allows parameters to be shared even within same CPD
- 3. Bayesian estimation
- 4. Hierarchical Bayes Framework
 - A softer version: parameters encouraged to be similar

Ex: Global Parameter Sharing

- Student is taking two courses c_1, c_2 with difficulties D_1, D_2
Grade G_i in class c_i depends on Intelligence and difficulty

- Assume that grades share the same CPDs

Intelligent student receives “A” in easy class
i.e., difficulty summarizes course challenge

- How to formalize this assumption

- For all choices of g, d and i we have

$$P(G_1=g|D_1=d, I=i) = P(G_2=g|D_2=d, I=i)$$

- Assumption does not imply grades are same

- But probability of getting a particular grade is same if class has same difficulty

- This is a network induced by a plate model

- Shown next

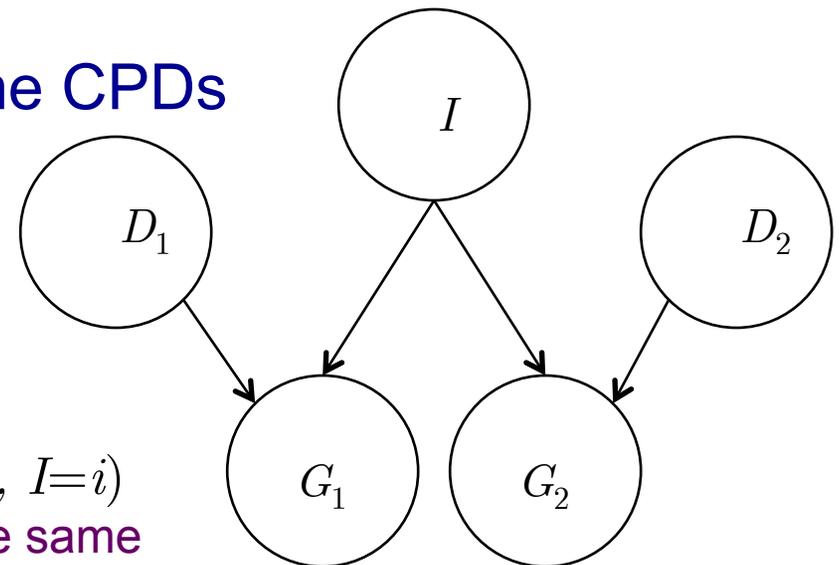
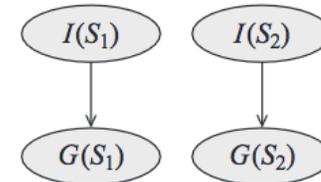
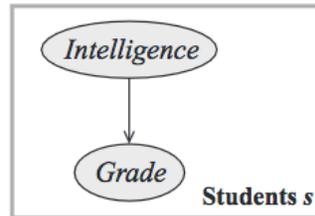


Plate models and induced BNs

Single Plate

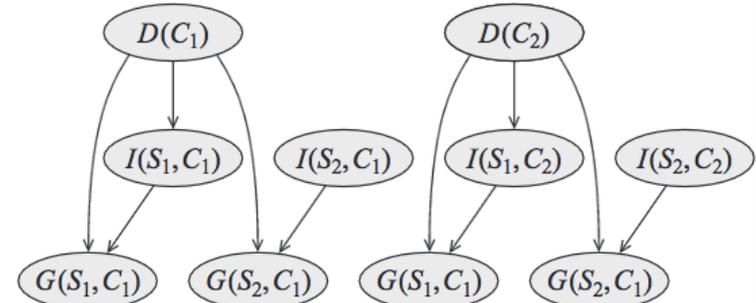
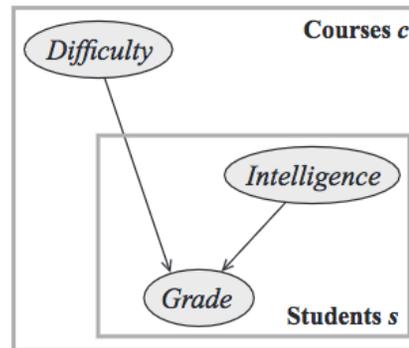
Multiple independent samples from same distribution



(a)

Nested Plates

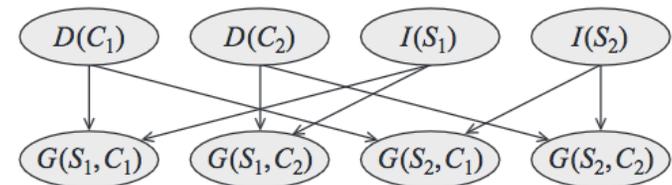
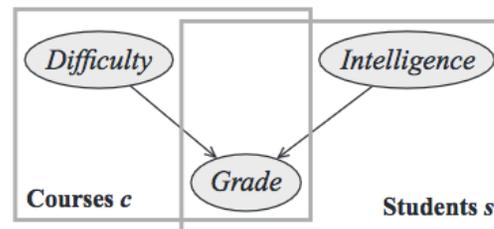
Multiple courses each with a separate set of students



(b)

Intersecting Plates

Multiple courses with overlapping sets of students



(c)

Plate models are a language for encoding models with repeated structure₅ and shared parameters

Likelihood with Global Shared Parameters

- Key to parameter estimation is in understanding likelihood function structure
- Consider network structure \mathcal{G} over a set of variables $\chi = \{X_1, \dots, X_n\}$ parameterized by θ
 - Variables X_i is associated with CPD $P(X_i | U_i, \theta)$
- Rather than assuming that each CPD has its own parameterization $\theta_{X_i | U_i}$ assume we have shared parameters used by multiple variables
- Thus sharing of parameters is global over \mathcal{G}

Partitioning of Parameters

- Assume that θ is partitioned into disjoint subsets $\theta^1, \dots, \theta^k$.
- With each subset associate a set of variables V^k so that V^1, \dots, V^k is disjoint partition of χ
 - For $X_i \in V^k$ the CPD of X_i depends only on θ^k ,
 - i.e., $P(X_i | \mathbf{U}_i, \theta) = P(X_i | \mathbf{U}_i, \theta^k)$
- We assume CPD is the same for all $X_i, X_j \in V^k$,
 - i.e., $P(X_i | \mathbf{U}_i, \theta) = P(X_j | \mathbf{U}_j, \theta^k)$
 - Makes sense only if $Val(X_i) = Val(X_j)$ & $Val(\mathbf{U}_i) = Val(\mathbf{U}_j)$
 - To avoid ambiguity, for any variable $X_i \in V^k$, we use y_k^l to range over values of X_i and w_k^l to range over possible values of parents

Max Likelihood of Shared Params

- Using the partitioning of variables, probability distribution is decomposed as

$$\begin{aligned}
 P(X_1, \dots, X_n | \theta) &= \prod_{i=1}^n P(X_i | Pa_{X_i}) \\
 &= \prod_{k=1}^K \prod_{X_i \in V_k} P(X_i | Pa_{X_i}, \theta) \quad \text{since } V^1, \dots, V^K \text{ defines a partition of } \chi \\
 &= \prod_{k=1}^K \prod_{X_i \in V_k} P(X_i | Pa_{X_i}, \theta^k) \quad \text{since } P(X_i | U_i, \theta) = P(X_i | U_i, \theta^k)
 \end{aligned}$$

- Let D be some assignment of values to

$$L(\theta : D) = \prod_{k=1}^K \prod_{X_i \in V_k} P(x_i | u_i, \theta^k)$$

which we can now rewrite

$$\begin{aligned}
 L(\theta, D) &= \prod_{k=1}^K \prod_{y_k, w_k} \prod_{\substack{X_i \in V^k \\ x_i = y_k, u_i = w_k}} \theta^k_{y_k w_k} \\
 &= \prod_{k=1}^K \prod_{y_k, w_k} \left(\theta^k_{y_k w_k} \right)^{\tilde{M}_k[y_k, w_k]} \quad \text{where we have a definition of counts } \tilde{M}_k[y_k, w_k] = \sum_{X_i \in V_k} I\{x_i = y_k, u_i = w_k\}
 \end{aligned}$$

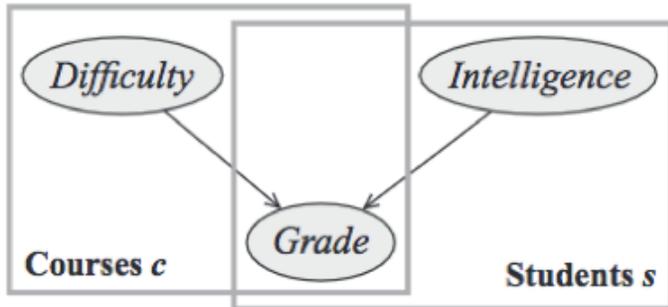
Thus we combine sufficient statistics from multiple variables across the network

- Thus maximum likelihood for each set of shared parameters is

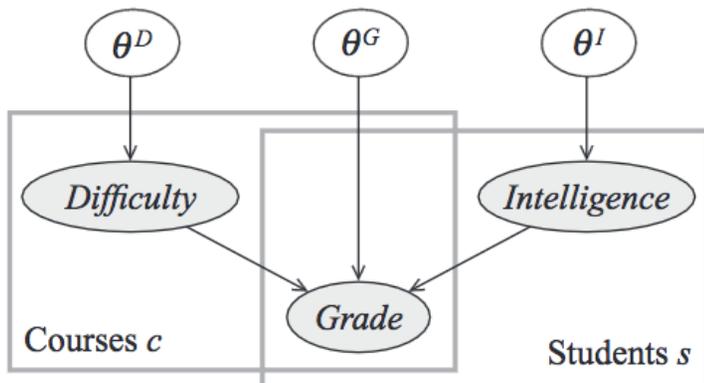
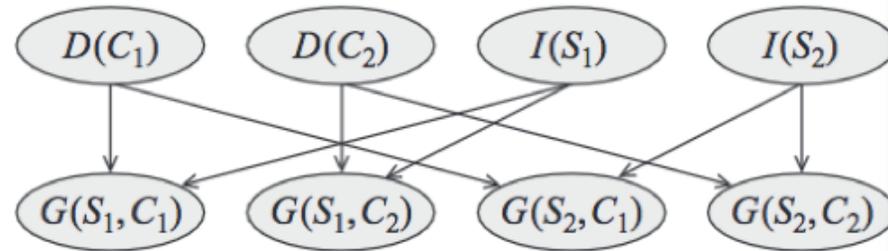
$$\hat{\theta}^k_{y_k w_k} = \frac{\tilde{M}_k[y_k, w_k]}{\tilde{M}_k[w_k]} \quad \text{Same estimate as with independent parameters using aggregate sufficient statistics}$$

Parameters as variables

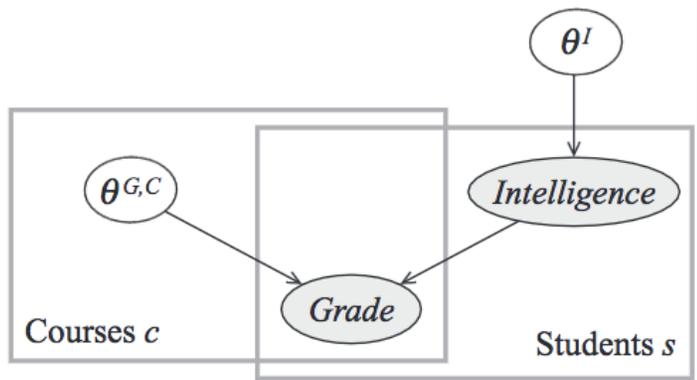
Plate Model



Ground network



Model where all parameters are global
 Same parameters are used in every CPD
 $P(G(s,c)|I(s),D(c))$ in the ground network



Each course has its own model for how difficulty affects grades.
 Parameters $\theta^{G,c}$ are used in CPDs of ground variables $G(s,c)$ for different values of s
 But there is no sharing between $G(s,c)$ and $G(s,c')$