

# CS 583: PROBABILISTIC GRAPHICAL MODELS

## TOPIC: TEMPLATE-BASED REPRESENTATIONS



**Mustafa Bilgic**



<http://www.cs.iit.edu/~mbilgic>



<https://twitter.com/bilgicm>

# MOTIVATION

- We have so far assumed that the set of variables  $\mathcal{X}$  is fixed
- In some domains, the size of  $\mathcal{X}$  can depend on the input
- Examples
  - Text, video, social networks, ...
- How can we specify the structure and the parameters for  $\mathcal{X}$ , where  $\mathcal{X}$  depends on the input?

# TEMPLATE-BASED MODELS

- We would like to specify a template model (both structure and parameters) where
  - The structure can be built automatically by using the template on the input
  - The parameters can be reused *shared*

# TEMPLATE-BASED MODEL EXAMPLES

- Temporal models
- Plate models
- Relational models

# TEMPORAL MODELS

- The state of the world, represented by a set of variables, evolves over time
  - $\mathcal{X}^{(0)}, \mathcal{X}^{(1)}, \dots, \mathcal{X}^{(T)}$
- Each possible world is a trajectory: an assignment to  $\mathcal{X}^{(0:T)}$
- Examples
  - Dynamic Bayesian networks
  - State observation models
  - Hidden Markov models
  - Kalman filters

# TEMPORAL MODELS

$$P\left(\mathcal{X}^{(0:T)}\right) = P\left(\mathcal{X}^{(0)}\right) \prod_{t=0}^{T-1} P\left(\mathcal{X}^{(t+1)} \mid \mathcal{X}^{(0:t)}\right)$$

Nothing magic; just the chain rule.

# MARKOV ASSUMPTION

$$x^{(t+1)} \perp x^{(0:t-1)} \mid x^{(t)}$$

Is this a reasonable assumption?

$x^{(0)}, x^{(1)}, x^{(2)}, \dots, x^{(t)}, \dots, x^{(\tau)}$

↑

# TEMPORAL MODELS

$$P\left(\mathcal{X}^{(0:T)}\right) = P\left(\mathcal{X}^{(0)}\right) \prod_{t=0}^{T-1} P\left(\mathcal{X}^{(t+1)} \mid \mathcal{X}^{(t)}\right)$$

$$X^{(0)} \rightarrow X^{(1)} \rightarrow X^{(2)} \rightarrow \dots \rightarrow X^{(T)}$$



# STATIONARY

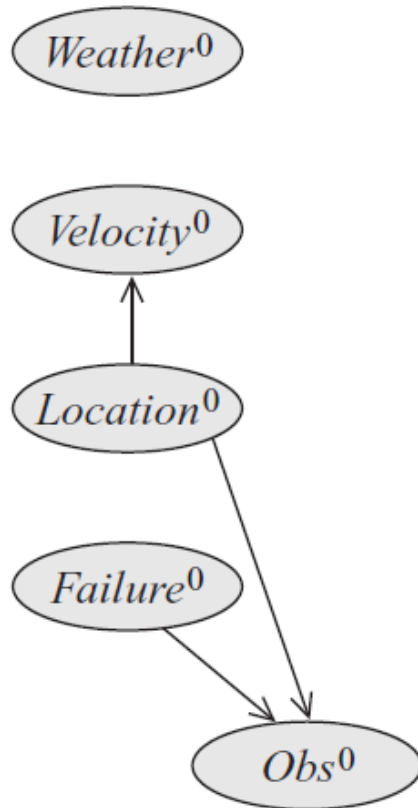
- We still need to specify the CPD for all variables in  $\mathcal{X}^{(0:T)}$
- A Markovian dynamic system is *stationary* (*time-invariant*, *homogeneous*) if  $P(\mathcal{X}^{(t+1)} \mid \mathcal{X}^{(t)})$  is same for all  $t$ . Then, the process can be represented using a *transition model*:

$$P(\mathcal{X}^{(t+1)} = \xi' \mid \mathcal{X}^{(t)} = \xi) = P(\mathcal{X}' = \xi' \mid \mathcal{X} = \xi)$$

# DYNAMIC BAYESIAN NETWORKS

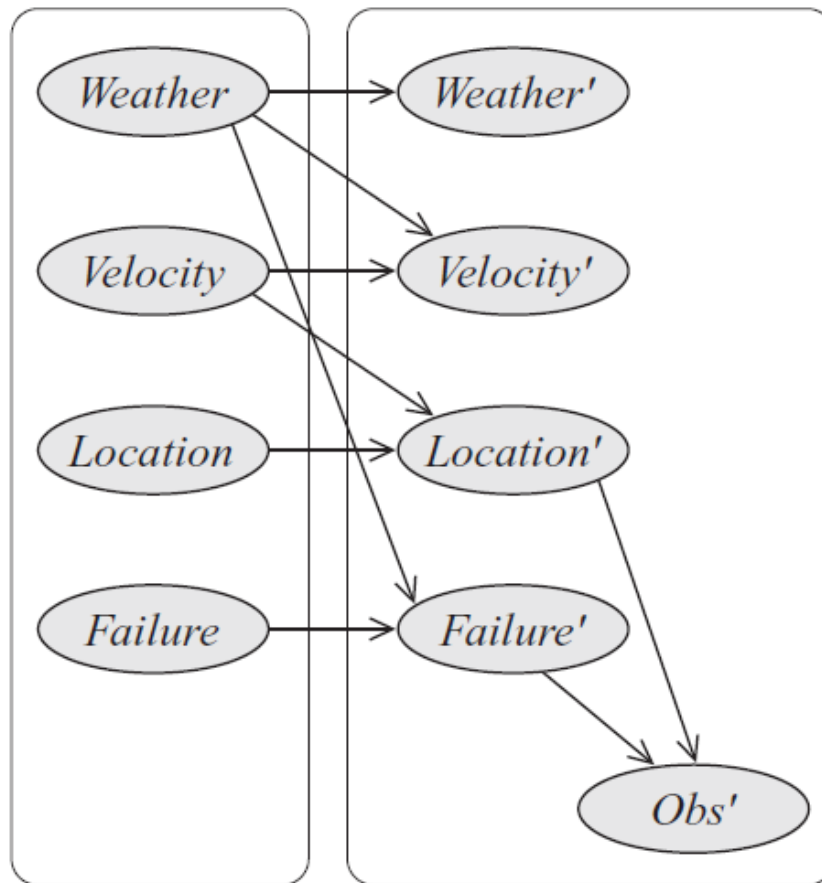
- Definition: A dynamic Bayesian network (DBN) is a pair  $\langle \mathcal{B}_0, \mathcal{B}_{\rightarrow} \rangle$ , where  $\mathcal{B}_0$  is a Bayesian network over  $\mathcal{X}^{(0)}$ , representing the initial distribution over states, and  $\mathcal{B}_{\rightarrow}$  is a 2-time-slice Bayesian network for the process.

# DBN EXAMPLE



Time slice 0

$\mathcal{B}_0$

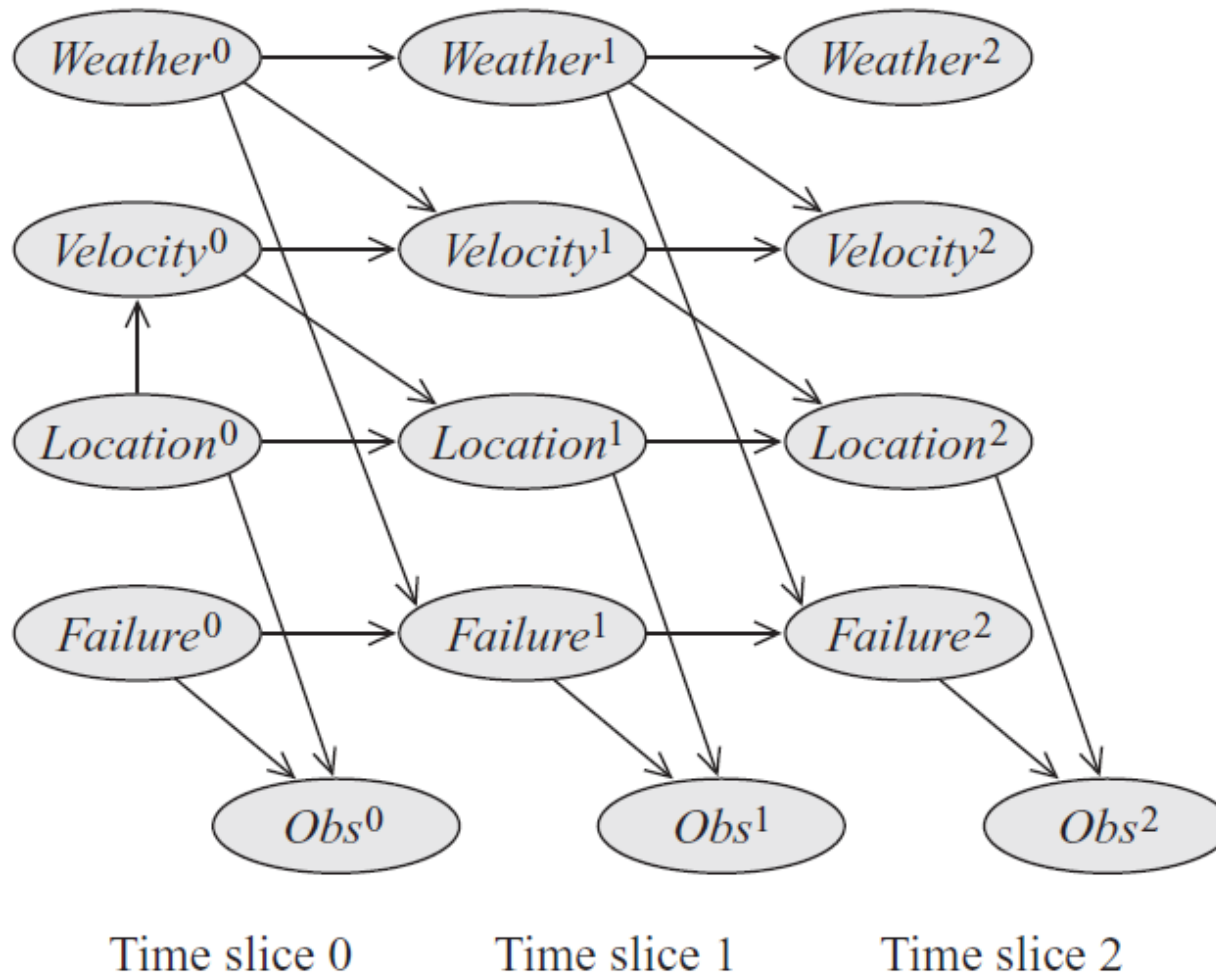


Time slice  $t$

Time slice  $t+1$

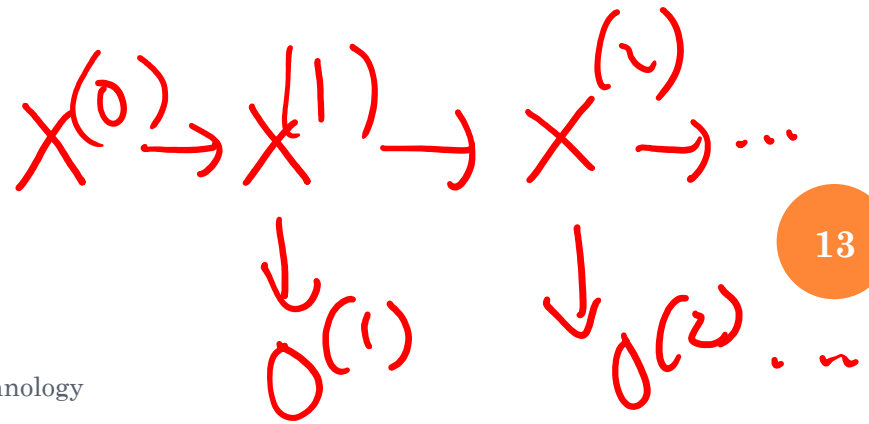
$\mathcal{B}_{\rightarrow}$

# DBN EXAMPLE



# STATE OBSERVATION MODELS

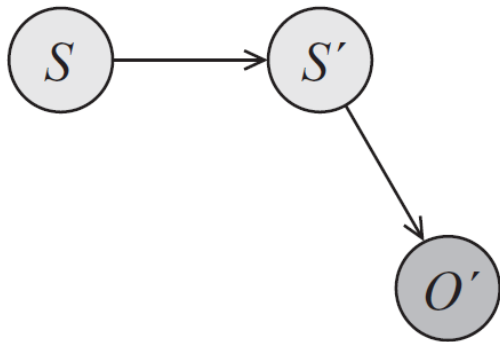
- System evolves on its own, with our observations occurring in a separate process
  - Our observations do not change the system
- Observations are obtained from a noisy sensor
- Two key assumptions
  - States are Markovian:  $\mathbf{X}^{(t+1)} \perp \mathbf{X}^{(0:t-1)} \mid \mathbf{X}^{(t)}$
  - Observations at  $t$ ,  $\mathbf{O}^{(t)}$ , are conditionally independent of the entire state sequence given the state at  $t$ ,  $\mathbf{X}^{(t)}$ :
    - $\mathbf{O}^{(t)} \perp \mathbf{X}^{(0:t-1)}, \mathbf{X}^{(t+1:\infty)} \mid \mathbf{X}^{(t)}$



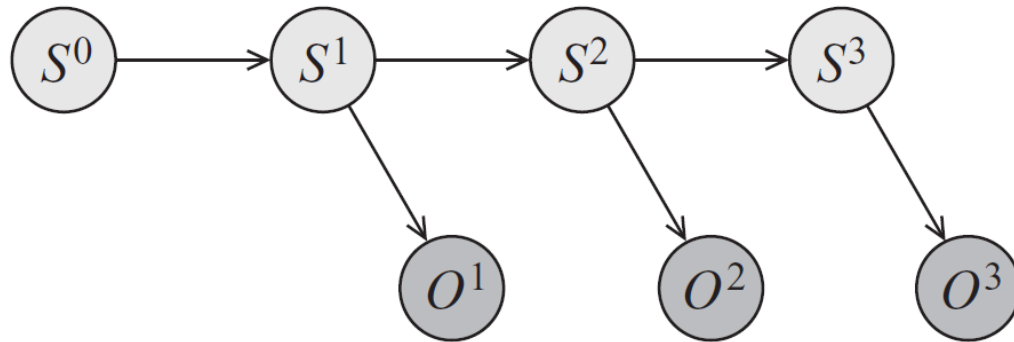
# HIDDEN MARKOV MODELS

$$P(S^u)$$

$$P(S^{t+1} | S^t)$$



(a)



(b)

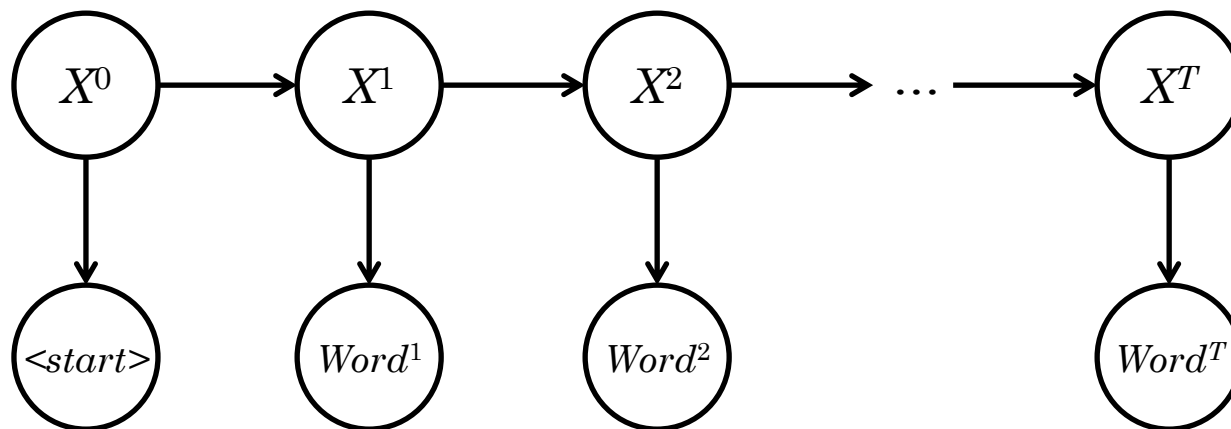
$$P(O^t | S^t)$$

# PART-OF-SPEECH TAGGING

- Task: Given a sentence, identify, for each word, the part of speech
  - E.g., noun, verb, adjective, adverb
- Part-of-speech depends on the context
  - E.g., “can”: modal or noun?
- Nearby words provide essential information
  - N + “can”: more likely to be a modal
  - “The” + “can”: more likely to be a noun

# PART-OF-SPEECH TAGGING USING HMMS

- States ( $\mathbf{X}$ ): Part-of-speech roles of the words
- Observations ( $\mathbf{O}$ ): The words themselves





# A GREAT READ

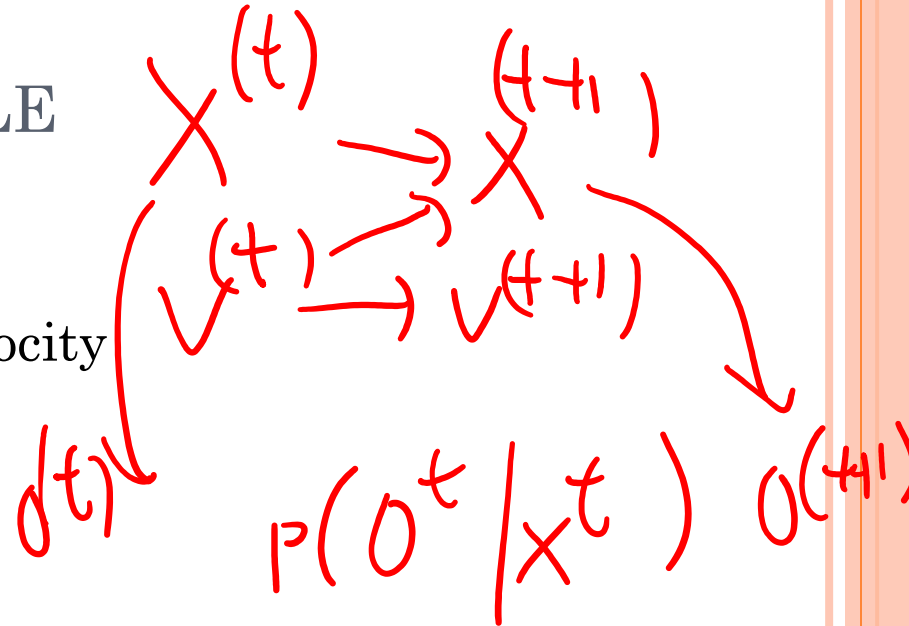
- “Speech and Language Processing” book by Jurafsky and Martin
- The full book is freely available at <https://web.stanford.edu/~jurafsky/slp3/>
- Chapter 8: Sequence Labeling for Parts of Speech and Named Entities
  - <https://web.stanford.edu/~jurafsky/slp3/8.pdf>
  - Introduces and discusses HMMs and CRFs

# KALMAN FILTERS

- A linear dynamical system
- Real-valued variables evolve linearly over time, with some Gaussian noise
- Can be considered a specific kind of dynamic Bayesian network where
  - All variables are continuous
  - All dependencies are linear Gaussian
- Often used for tracking systems, for example tracking airplanes using radar data

# KALMAN FILTER EXAMPLE

- Tracking a vehicle's position
- X: current position, V: current velocity
- $P(X'|X, V) = X + V\Delta + \mathcal{N}(0; \sigma_X^2)$
- $P(V'|V) = V + \mathcal{N}(0; \sigma_V^2)$
- Often, there is also a noisy observation; for e.g., the GPS signal as a noisy measurement of X.
- How would you model the above system using a dynamic Bayesian network?



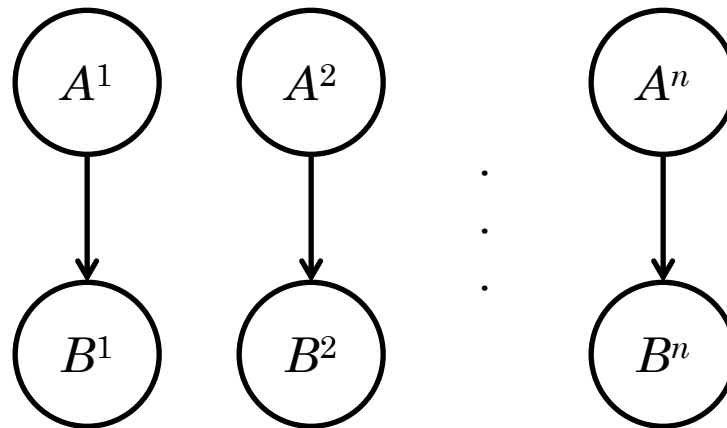
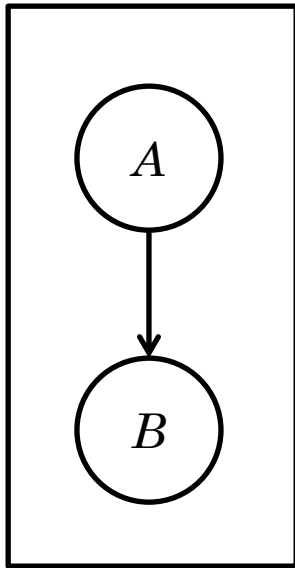
# PLATE MODELS

- **Motivation:** Assume we have the following simple network
  - $D \rightarrow G \leftarrow I$
- There are multiple students in a given class, say CS101
- Given the grade of John, we would like to reason about the grade distribution of other students
  - For e.g., if John got an A, we would like to be able to reason that CS101 might be an easy class and thus other students in CS101 have a higher chance of getting As as well

# PLATE MODELS - INFORMALLY

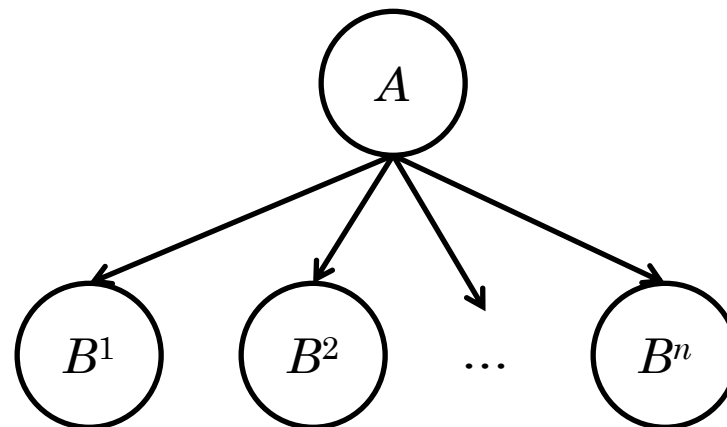
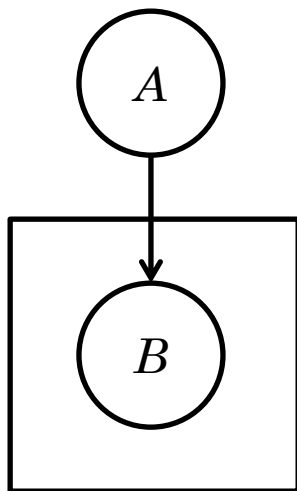
- A plate is a rectangle with a Bayesian network inside
- The rectangle tells us to replicate the same structure *and* the parameters multiple times
- Let's see some examples

## EXAMPLE – ALL INSIDE



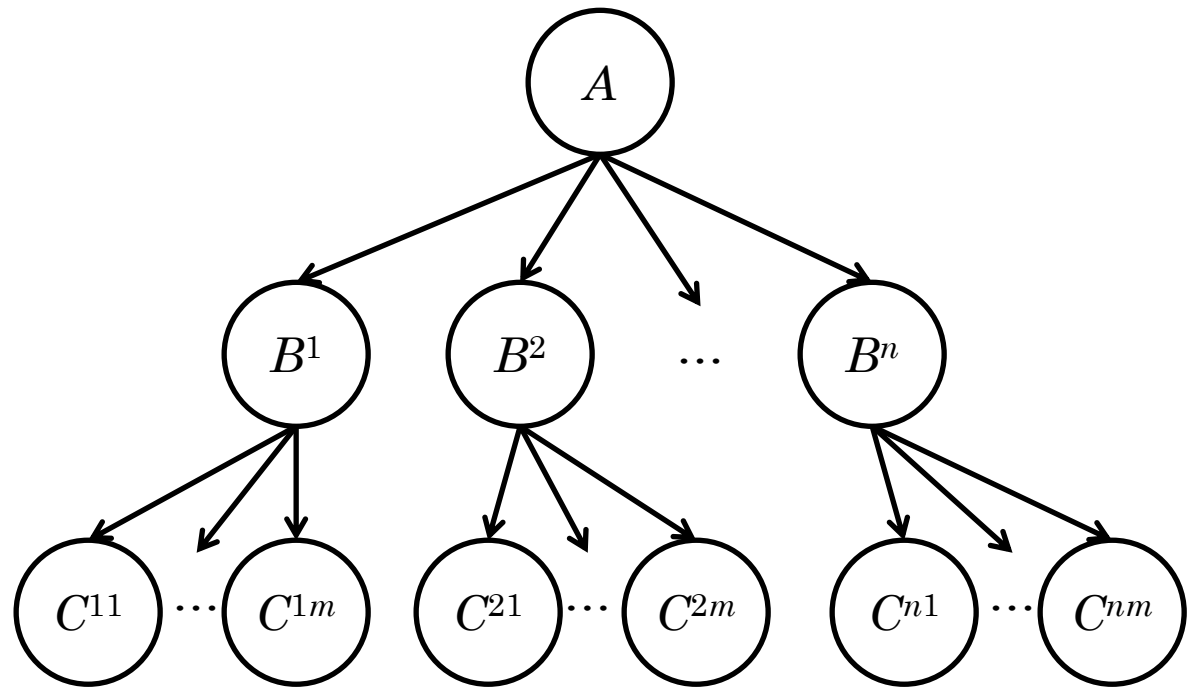
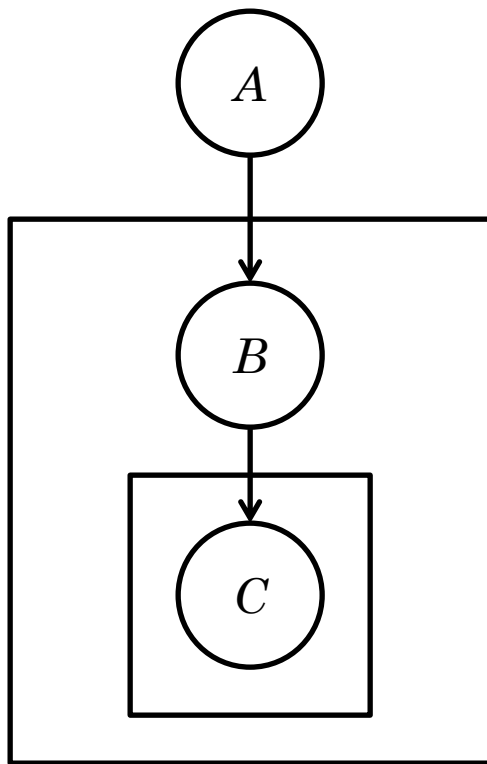
The marginal  $P(A)$  and the conditional  $P(B|A)$  are copied for each replica.

## EXAMPLE – SOME OUTSIDE



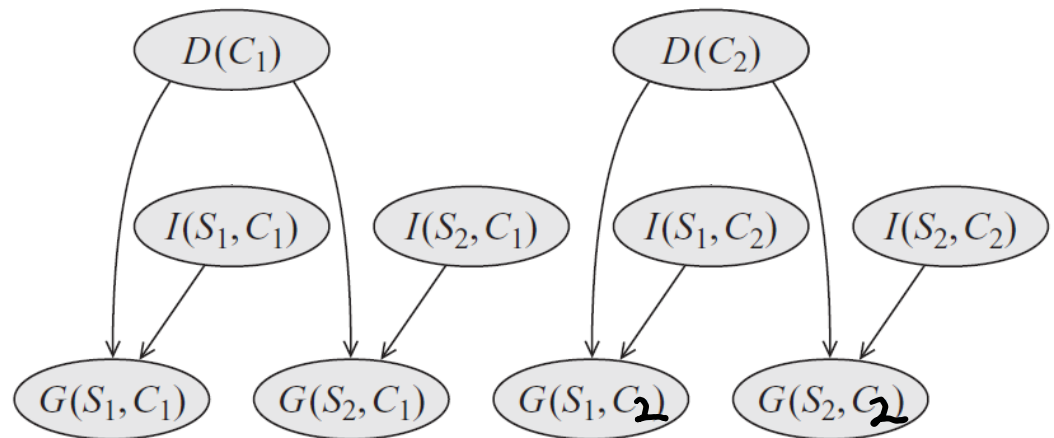
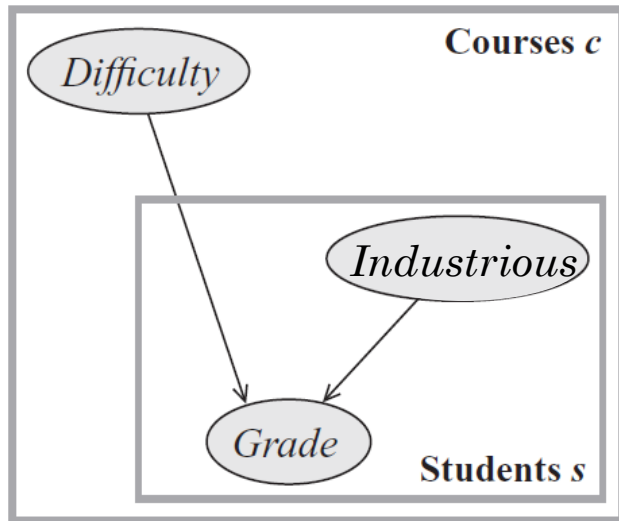
The conditional  $P(B|A)$  is copied for each replica.

## EXAMPLE – NESTED



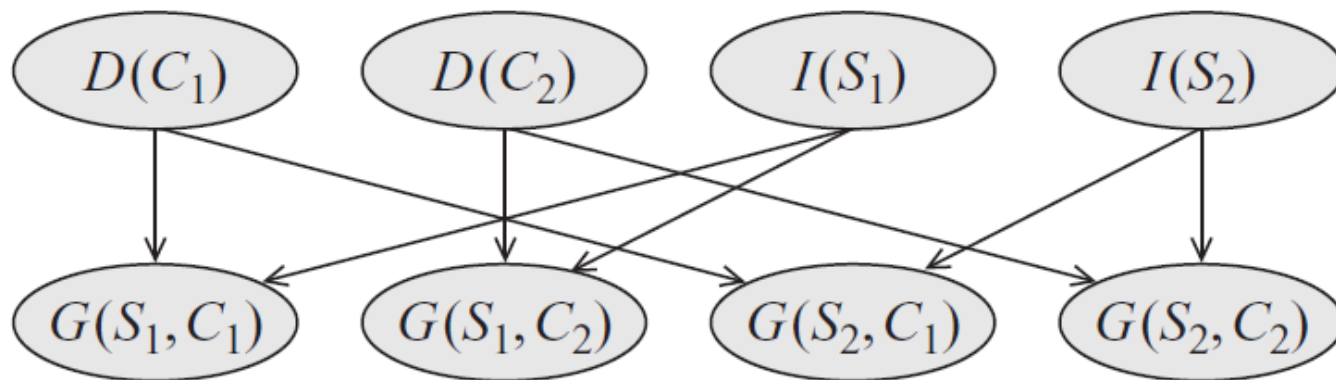
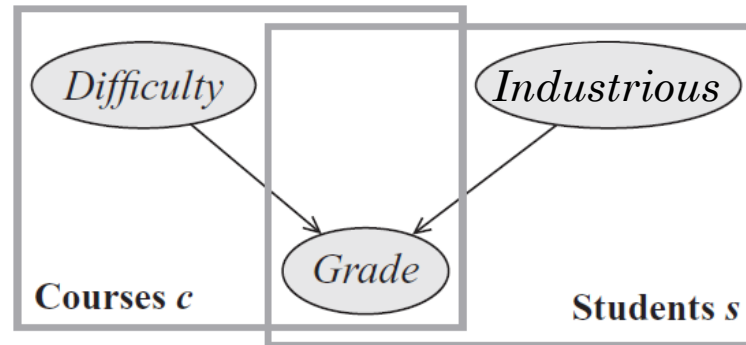


# MULTIPLE COURSES WITH EACH ITS OWN SET OF STUDENTS

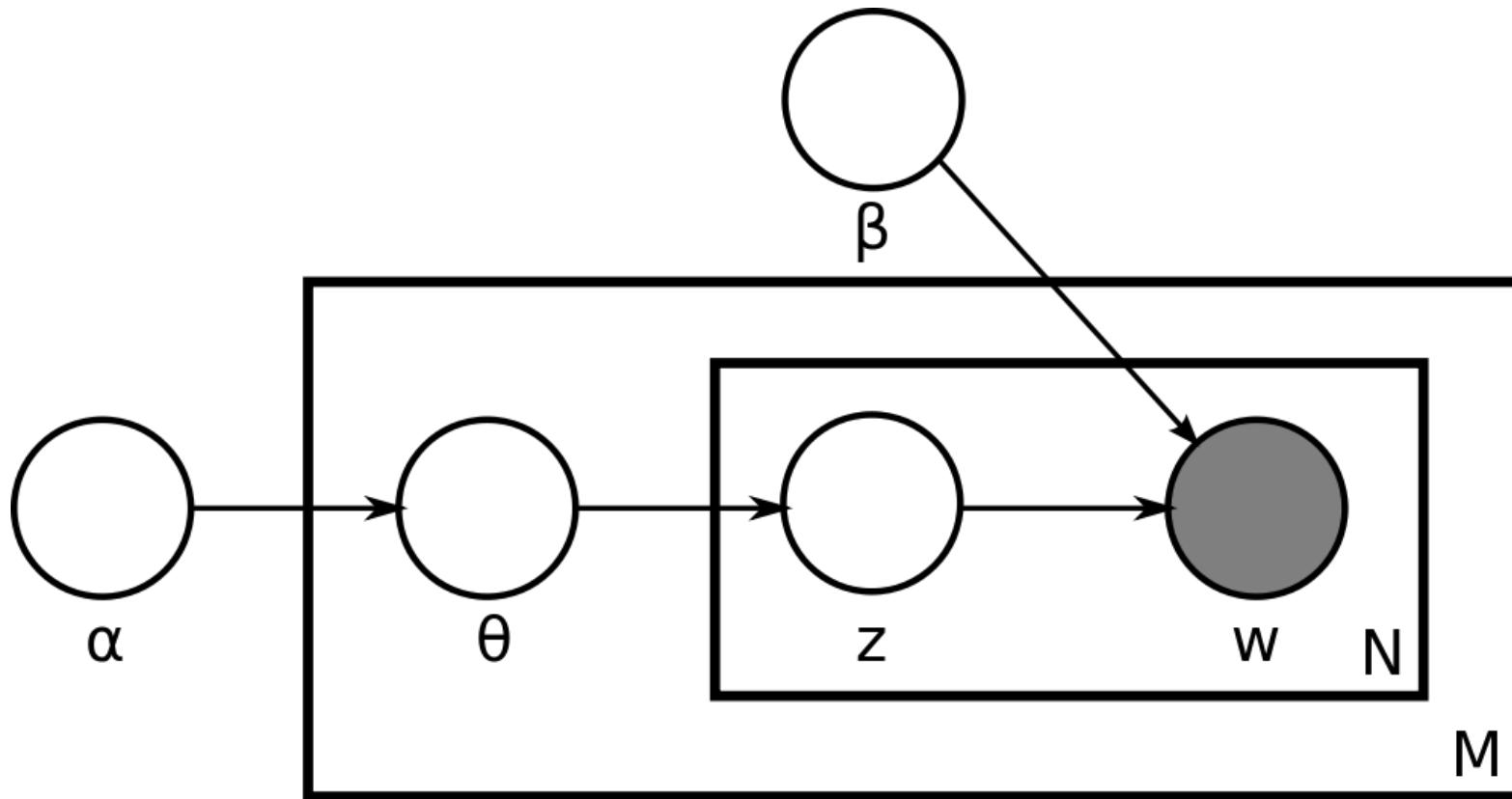


We can reason across students in the same course, but we cannot reason across courses through common students, because the student variables are replicated across courses.

# MULTIPLE COURSES AND MULTIPLE STUDENTS



# LATENT DIRICHLET ALLOCATION



<https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>

<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.LatentDirichletAllocation.html>

# LIMITATIONS OF PLATE MODELS

- The grade is replicated for all possible combinations of courses and students. What happens when a student takes some classes but not others?
- How can you represent a variable of one object depends on a variable of another object of the same type?
- How can you represent temporal models?

# PROBABILISTIC RELATIONAL MODELS

- Specify when a dependency holds
  - For e.g., the dependency between a course variable and a student variable is contingent upon whether the student registered that course
- Largely based on database concepts
  - Entities and relationships

# STATISTICAL RELATIONAL LEARNING (SRL)

- Last decade has seen tremendous amount of work on probabilistic relational modeling
- Examples
  - Probabilistic Relational Models (PRM)
  - Relational Markov Networks (RMN)
  - Inductive Logic Programming (ILP)
  - Conditional Random Fields (CRF)
  - Relational Dependency Networks (RDN)
  - Bayesian Logic Programming (BLP)
  - Markov Logic Networks (MLN)
- <http://www.cs.umd.edu/srl-book/>