

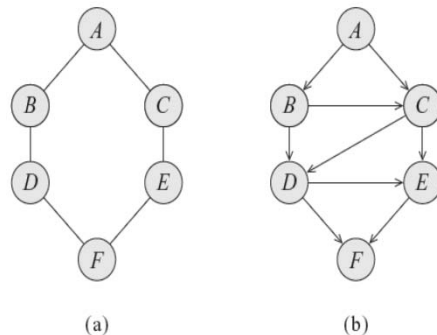
Lecture 7: February 4, 2015
cs 573: Probabilistic Reasoning
Professor Nevatia
Spring 2015

Review

- Assignment #2 due Feb 9
- Update on TA office hours: is one day per week enough?
- Last lecture:
 - Intro to Markov Networks (undirected graphs)
 - Factors and Factor product
 - Reduced Factors
 - Gibbs distribution
 - Local and Global Independences
 - BNs to MNs
- Today's objective
 - Chordal graphs
 - Log Linear Models
 - CRFs

$$MN \Rightarrow BN$$

- Much harder in this direction



- Start with nodes in MN (in some order). Construct minimal I-map BN; example from one order is shown above
 - Consider A,B,C,D,E,F order
 - A,B and A,C obviously connected but we also need edge from B to C as in MN, ($B \not\perp C|A$)
 - Consider D: in MN ($D \not\perp C|B$) so edge from C to D in BN..
 - from C to D in BN
- Note that directed graph is chordal (triangulated)
- In general, minimal I-map, $G(BN)$, for any MN, H , is necessarily chordal

Chordal Graph Properties

- If H is a non-chordal MN, there is no BN, G , that is a perfect map of H .
- If H is a chordal MN, there exists a BN, G , such that $I(H) = I(G)$
- Chordal graphs also map to *clique trees* that permit efficient inference algorithms (ch. 10); we will postpone definition of clique-trees until we make use of them in the inference algorithms
- Venn diagram
 - BN and MN subsets of P
 - Intersection of BN and MN is set of chordal graphs

Log Linear Model

- Rewrite $\Phi(D) = e^{-\epsilon(D)}$
- $\epsilon(D) = -\ln \Phi(D)$, ϵ is called an *energy* function (Φ is called clique *potential*)

$$P(X_1, \dots, X_n) \propto \exp \left[- \sum_{i=1}^m \epsilon_i(D_i) \right]$$

$\epsilon_1(A, B)$	$\epsilon_2(B, C)$	$\epsilon_3(C, D)$	$\epsilon_4(D, A)$
$a^0 \quad b^0 \quad -3.4$	$b^0 \quad c^0 \quad -4.61$	$c^0 \quad d^0 \quad 0$	$d^0 \quad a^0 \quad -4.61$
$a^0 \quad b^1 \quad -1.61$	$b^0 \quad c^1 \quad 0$	$c^0 \quad d^1 \quad -4.61$	$d^0 \quad a^1 \quad 0$
$a^1 \quad b^0 \quad 0$	$b^1 \quad c^0 \quad 0$	$c^1 \quad d^0 \quad -4.61$	$d^1 \quad a^0 \quad 0$
$a^1 \quad b^1 \quad -2.3$	$b^1 \quad c^1 \quad -4.61$	$c^1 \quad d^1 \quad 0$	$d^1 \quad a^1 \quad -4.61$
(a)	(b)	(c)	(d)
$\phi_1(A, B)$	$\phi_2(B, C)$	$\phi_3(C, D)$	$\phi_4(D, A)$
$a^0 \quad b^0 \quad 30$	$b^0 \quad c^0 \quad 100$	$c^0 \quad d^0 \quad 1$	$d^0 \quad a^0 \quad 100$
$a^0 \quad b^1 \quad 5$	$b^0 \quad c^1 \quad 1$	$c^0 \quad d^1 \quad 100$	$d^0 \quad a^1 \quad 1$
$a^1 \quad b^0 \quad 1$	$b^1 \quad c^0 \quad 1$	$c^1 \quad d^0 \quad 100$	$d^1 \quad a^0 \quad 1$
$a^1 \quad b^1 \quad 10$	$b^1 \quad c^1 \quad 100$	$c^1 \quad d^1 \quad 1$	$d^1 \quad a^1 \quad 100$
(a)	(b)	(c)	(d)

Log Linear Model

- Often, we will find that the logarithmic representation is easier to manipulate and numerically more stable
 - Adding rather than multiplying many numbers.
 - Minimizing energy is equivalent to maximizing probability.
- Feature Functions
 - Feature $f(\mathbf{D})$ is a function from $\text{val}(\mathbf{D})$ to \mathbb{R} .
 - *Indicator* feature: has value 1 for some value, y in $\text{val}(\mathbf{D})$ (*i.e.* for some assignment of values to variables in \mathbf{D})
 - Another feature: $\varepsilon(A_1, A_2) = -3$ if $A_1 = A_2$; 0 otherwise.
- $\varepsilon = \sum_k w_k f_k(\mathbf{D}_k)$; w_k weight parameters; features typically but not necessarily binary, multiple features over same clique acceptable
- $P = (1/Z) e^{-\sum_k w_k f_k(\mathbf{D}_k)}$

Feature Representation Example

- We can represent arbitrary factors using feature functions
- Consider $\epsilon_1(a,b)$ from misconception example

$\epsilon_1(A, B)$		
a^0	b^0	-3.4
a^0	b^1	-1.61
a^1	b^0	0
a^1	b^1	-2.3

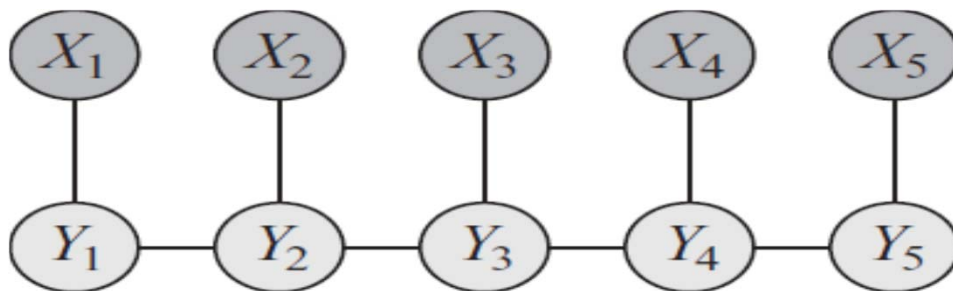
- Define $f_1(A, B) = 1$ when $A=0$ and $B=0$; zero otherwise
 - This provides first entry in the factor, set weight = -3.4
- $f_2(A, B) = 1$ when $A=0$ and $B=1$, zero otherwise
 - Provides second entry into the table, set weight = -1.61
- ...
- Flexible methodology: can assign weights to arbitrary functions of variables in a clique
 - Setting of these weights may not always be intuitive and learning them may be more complex than learning conditional probabilities

Overparameterization

- It is easy to overparameterize a Markov Network
 - Nodes can be in multiple cliques, evidence for them can be distributed differently
- Canonical parameterization
 - Factors for each clique (not just the maximal cliques)
 - We skip details (see section 4.4.2.1)
- If we use feature functions, some linear relations between them apply (i.e. they are not linearly independent)
 - Can solve for a minimal set of features
 - Again, we skip details, see section 4.4.2.2

Conditional Random Fields (CRFs)

- Encodes conditional distribution $P(\mathbf{Y}|\mathbf{X})$; Y is a set of target variables (unknowns), X is a set of observed variables. Note that MRF encodes $P(\mathbf{Y}, \mathbf{X})$.
- First an example:



A conditional random field is an undirected graph \mathcal{H} whose nodes correspond to $\mathbf{X} \cup \mathbf{Y}$; the network is annotated with a set of factors $\phi_1(\mathbf{D}_1), \dots, \phi_m(\mathbf{D}_m)$ such that each $\mathbf{D}_i \not\subseteq \mathbf{X}$. The network encodes a conditional distribution as follows:

$$\begin{aligned} P(\mathbf{Y} | \mathbf{X}) &= \frac{1}{Z(\mathbf{X})} \tilde{P}(\mathbf{Y}, \mathbf{X}) \\ \tilde{P}(\mathbf{Y}, \mathbf{X}) &= \prod_{i=1}^m \phi_i(\mathbf{D}_i) \\ Z(\mathbf{X}) &= \sum_{\mathbf{Y}} \tilde{P}(\mathbf{Y}, \mathbf{X}). \end{aligned} \tag{4.11}$$

Two variables in \mathcal{H} are connected by an (undirected) edge whenever they appear together in the scope of some factor.

CRFs: General Definition

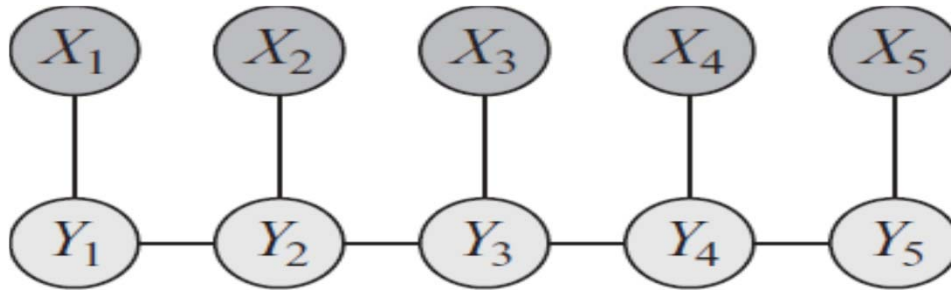
A conditional random field is an undirected graph \mathcal{H} whose nodes correspond to $\mathbf{X} \cup \mathbf{Y}$; the network is annotated with a set of factors $\phi_1(\mathbf{D}_1), \dots, \phi_m(\mathbf{D}_m)$ such that each $\mathbf{D}_i \not\subseteq \mathbf{X}$. The network encodes a conditional distribution as follows:

$$\begin{aligned} P(\mathbf{Y} \mid \mathbf{X}) &= \frac{1}{Z(\mathbf{X})} \tilde{P}(\mathbf{Y}, \mathbf{X}) \\ \tilde{P}(\mathbf{Y}, \mathbf{X}) &= \prod_{i=1}^m \phi_i(\mathbf{D}_i) \\ Z(\mathbf{X}) &= \sum_{\mathbf{Y}} \tilde{P}(\mathbf{Y}, \mathbf{X}). \end{aligned} \tag{4.11}$$

Two variables in \mathcal{H} are connected by an (undirected) edge whenever they appear together in the scope of some factor. ■

Note: connections do not have to be between adjacent nodes in a chain only, one target variable node may be connected to multiple evidence nodes. Flexible due to easy use of feature functions and log-linear models. Linear chain graphs are natural for many problems, e.g. text, activity analysis, robot navigation....

Conditional Random Fields (CRFs)



$$P(\mathbf{Y} \mid \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \tilde{P}(\mathbf{Y}, \mathbf{X})$$

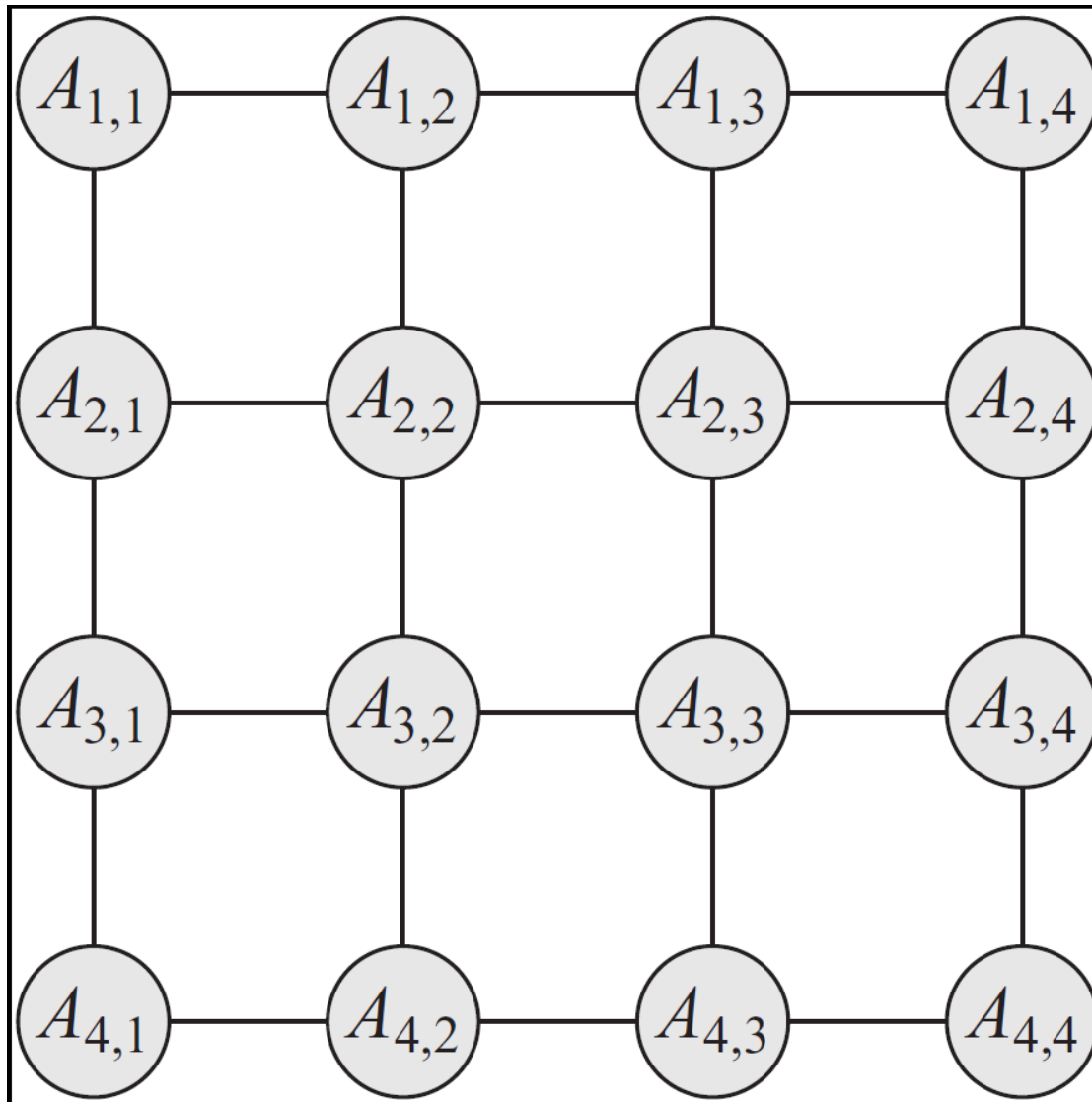
$$\tilde{P}(\mathbf{Y}, \mathbf{X}) = \prod_{i=1}^{k-1} \phi(Y_i, Y_{i+1}) \prod_{i=1}^k \phi(Y_i, X_i)$$

$$Z(\mathbf{X}) = \sum_{\mathbf{Y}} \tilde{P}(\mathbf{Y}, \mathbf{X}).$$

Note: no factors over just variables in \mathbf{X} ; Z is now a function of \mathbf{X} , the values of the observed variables. Connections do not have to be between adjacent nodes in a chain only, one target variable node may be connected to multiple evidence nodes.

This model can not predict $P(\mathbf{X})$ or $P(\mathbf{Y}, \mathbf{X})$. In this example, functions are same across the network.

Pairwise Markov Random Field (MRF)



Notes:

- a) could represent as a single clique, but would be highly complex
- b) Pairwise MRFs can not represent all distributions
- c) Commonly used in vision and other applications for simplicity
- d) CRF if each node is also connected to an observable

MRFs/CRFs in Computer Vision (Illustrations only)

- Consider each node also connected to another “observation” node (node whose value is given); this makes MRF into a CRF (Conditional Random Field) but term MRF is still commonly used.
- *Denoising*: Observed data is corrupted by noise. Goal is to label nodes with the *correct* numbers. Unary function has high value if label is close to observed value. Binary factor prefers similar labels for neighbors. MAP solution gives the denoised estimate.
- *Segmentation*: Label each pixel as “foreground” or “background”. Unary function provides a label according to some visual features at the point (color, texture, brightness...). Binary function prefers continuity (*i.e.* neighbors with different labels have low values). MAP solution gives best segmentation, given this model. Results are considerably better by considering neighbors.
- Labeling: assign class labels to pixels/regions

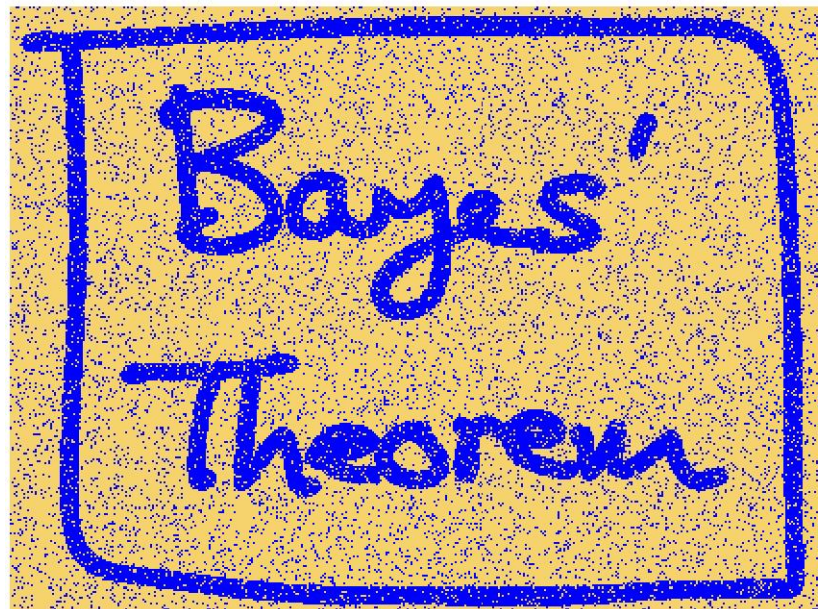
Ising Model

- Originally developed for statistical physics (modeling influence of neighboring atoms)
- Let random variables be binary (have +1 or -1 values)
- $\varepsilon_{i,j} (x_i, x_j) = w_{i,j} x_i x_j$
- Positive contribution of $w_{i,j}$ when $x_i = x_j$; negative $w_{i,j}$ otherwise.
- Example: same spin for two atoms in ferromagnetic analysis; same intensity value or label in image analysis ...
- Combine with unary energy term (based on probability of each random variable value by itself).
- Efficient algorithms exist for computing MAP solution for this model
- Boltzmann machine: variables take values $\{0,1\}$
 - Used to model a simple neuron behavior
- *Potts* model is a generalization for multiple valued variables

Illustration: Image De-Noising



Original
Image

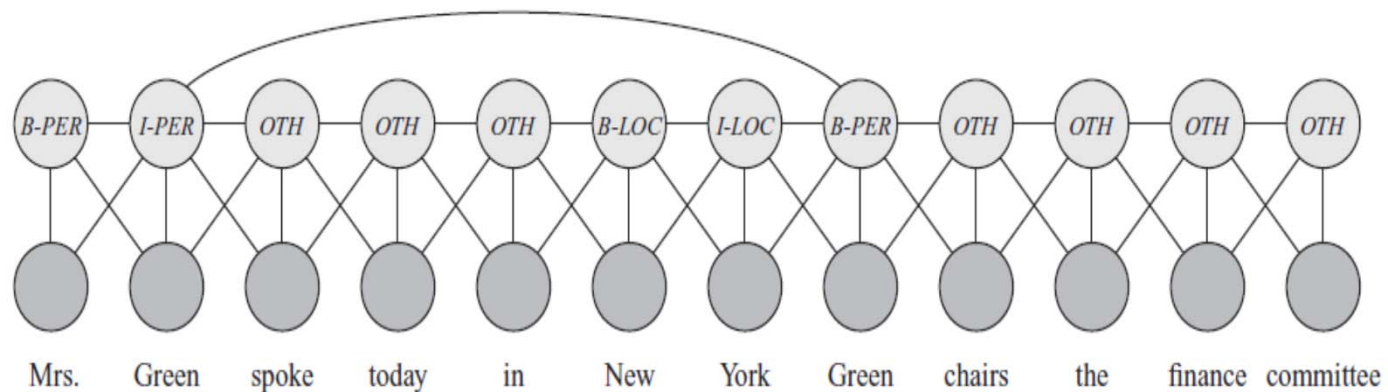


Noisy Image

From: Bishop Book

CRF: Text Analysis Example

- Network for text analysis. CRFs were originally introduced for text analysis (~2001) though now they find common in use in many applications. Note multiple connections to observation nodes. Also, model is not quite a linear chain (contrasts with HMMs which we haven't studied yet).



KEY

<i>B-PER</i>	Begin person name	<i>I-LOC</i>	Within location name
<i>I-PER</i>	Within person name	<i>OTH</i>	Not an entity
<i>B-LOC</i>	Begin location name		

(a)

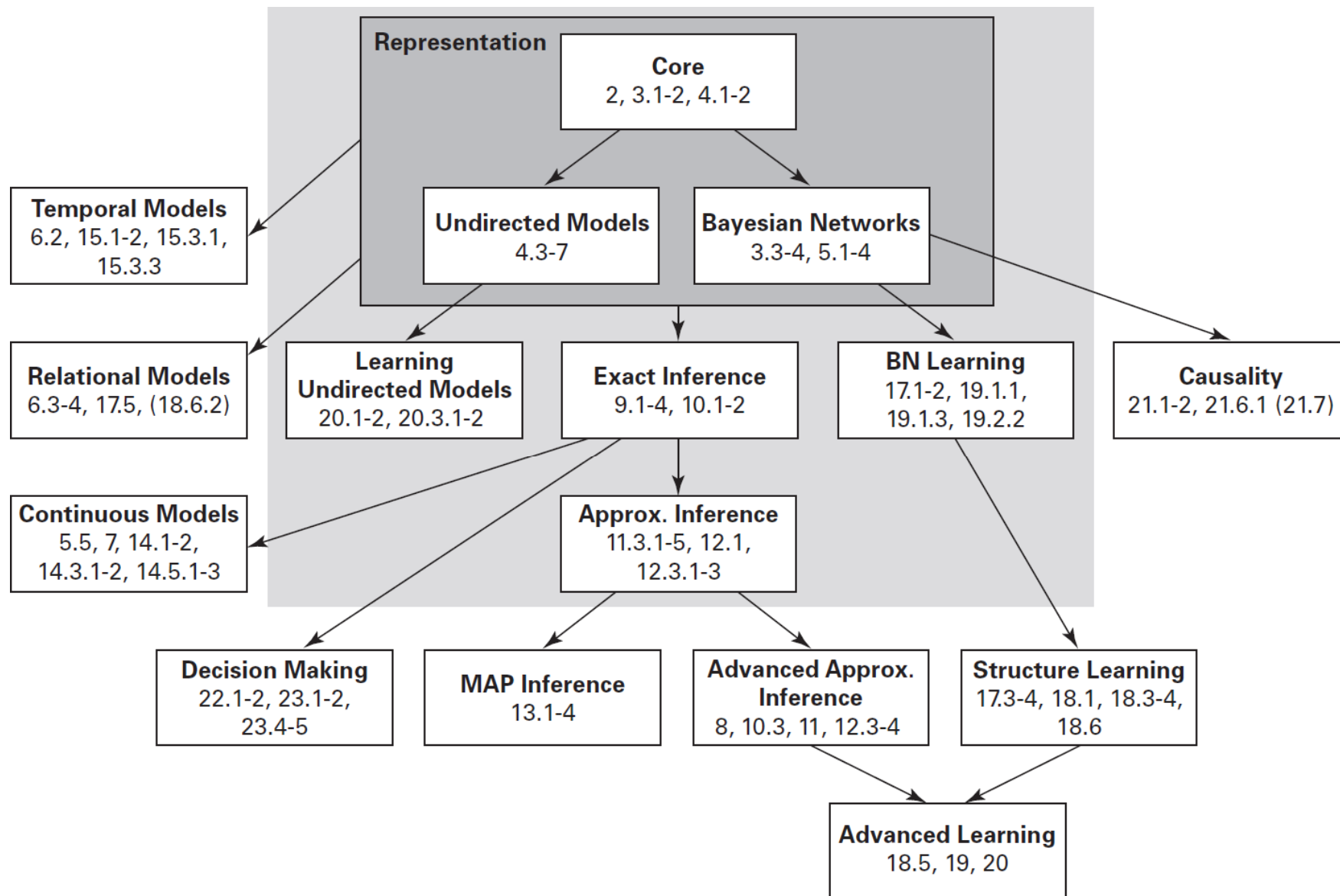
Summary of Topics Covered so Far

- Foundations (chapter 2) basic concepts of probability theory
- Bayesian networks and Markov networks
 - Basic representation
 - Factorization, conditional independences
 - I-maps, P-maps
 - Graphs \Leftrightarrow distributions,
 - BN \Leftrightarrow MN
- Local models (just summarized)
- We have covered chapters 2-5 except for
 - Sections 3.4.3, 4.4.2, 4.6.2, 5.3.2, 5.4.3, 5.4.4, 5.5.1
 - We have (mostly) skipped proofs of theorems

What More is There (in Representation)?

- Temporal models (ch. 6)
 - Future state depends on current state, dependency is not time dependent
 - Can describe relations by a “template” rather than a long sequence of variables
 - We will cover representation together with temporal inference later in the course
- “Plate” Models (ch. 6)
 - Describe multiple variables with same parameters, *i.e.* tossing hundred fair coins; will study along with parameter learning (time permitting)
- Gaussian Networks (ch. 7): Useful in modeling continuous variables: will study a bit later
- Exponential Family: General framework that includes many important distributions classes (Binomial, multinomial, Gaussian...)
 - Will likely not be covered in the class

Book Plan



Next Class

- Read sections 9.1 to 9.3 of the KF book