

Recap of Basic Prob. Concepts

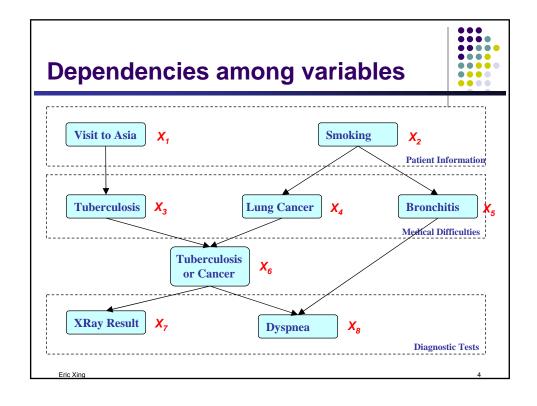


 Representation: what is the joint probability dist. on multiple variables?

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8,)$$

- How many state configurations in total? --- 28
- Are they all needed to be represented?
- Do we get any scientific/medical insight?
- Learning: where do we get all this probabilities?
 - Maximal-likelihood estimation? but how many data do we need?
 - Where do we put domain knowledge in terms of plausible relationships between variables, and plausible values of the probabilities?
- Inference: If not all variables are observable, how to compute the conditional distribution of latent variables given evidence?

Eric Xino



Probabilistic Graphical Models



- · Represent dependency structure with a graph
 - Node <-> random variable
 - Edges encode dependencies
 - Absence of edge -> conditional independence
 - Directed and undirected versions



- Why is this useful?
 - A language for communication
 - A language for computation
 - A language for development
- Origins:
 - Wright 1920's
 - Independently developed by Spiegelhalter and Lauritzen in statistics and Pearl in computer science in the late 1980's

Eric Xing

5

Probabilistic Graphical Models, con'd



□ If X_i 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,



$$\begin{split} & P(X_p \ X_2 \ X_3 \ X_4 \ X_5 \ X_6 \ X_7 \ X_8) \\ & = P(X_1) \ P(X_2) \ P(X_3 | X_1) \ P(X_4 | X_2) \ P(X_5 | X_2) \\ & P(X_6 | X_2 \ X_4) \ P(X_7 | X_6) \ P(X_8 | X_2 \ X_6) \end{split}$$

- Why we may favor a PGM?
 - Representation cost: how many probability statements are needed?

2+2+4+4+4+8+4+8=36, an 8-fold reduction from 28!

- Algorithms for systematic and efficient inference/learning computation
 - Exploring the graph structure and probabilistic (e.g., Bayesian, Markovian) semantics
- Incorporation of domain knowledge and causal (logical) structures

Eric Xing

Two types of GMs



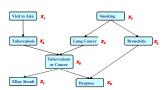
- Directed edges give causality relationships (Bayesian Network or Directed Graphical Model):
- Undirected edges simply give (physical or symmetric) correlations between variables (Markov Random Field or Undirected Graphical model):

Eric Xing

7

Bayesian Network: Factorization Theorem





$$\begin{split} & P(X_p \ X_2 \ X_3 \ X_{\phi} \ X_5 \ X_{\phi} \ X_7 \ X_8) \\ & = P(X_1) \ P(X_2) \ P(X_3|\ X_1) \ P(X_4|\ X_2) \ P(X_5|\ X_2) \\ & P(X_6|\ X_3 \ X_4) \ P(X_7|\ X_6) \ P(X_8|\ X_5 \ X_6) \end{split}$$

• Theorem:

Given a DAG, The most general form of the probability distribution that is consistent with the graph factors according to "node given its parents":

$$P(\mathbf{X}) = \prod_{i=1}^{d} P(X_i \mid \mathbf{X}_{\pi_i})$$

where X_{π_i} is the set of parents of x_i , d is the number of nodes (variables) in the graph.

Eric Xing

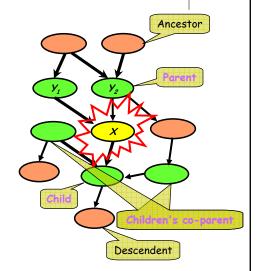
Bayesian Network: Conditional Independence Semantics



Structure: DAG

- Meaning: a node is conditionally independent of every other node in the network outside its Markov blanket
- Local conditional distributions (CPD) and the DAG completely determine the joint dist.
- Give causality relationships, and facilitate a generative process

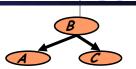
Eric Xino



Local Structures & Independencies



Fixing B decouples A and C
 "given the level of gene B, the levels of A and C are independent"



Cascade

Knowing B decouples A and C
 "given the level of gene B, the level gene A provides no extra prediction value for the level of gene C"



V-structure

Knowing C couples A and B
 because A can "explain away" B w.r.t. C

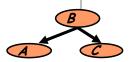
"If A correlates to C, then chance for B to also correlate to B will decrease"

• The language is compact, the concepts are rich!

Eric Xino

A simple justification





Eric Xing

11

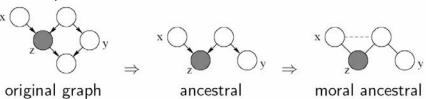
Graph separation criterion



• D-separation criterion for Bayesian networks (D for Directed edges):

Definition: variables x and y are *D-separated* (conditionally independent) given z if they are separated in the *moralized* ancestral graph

• Example:

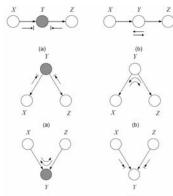


Eric Xino

Global Markov properties of DAGs



X is d-separated (directed-separated) from Z given Y if we can't send a ball from any node in X to any node in Z using the "Bayes-ball" algorithm illustrated bellow (and plus some boundary conditions):



 Defn: *I(G)*=all independence properties that correspond to dseparation:

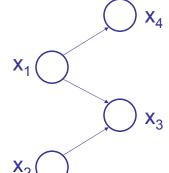
$$I(G) = \{X \perp Z | Y : dsep_G(X; Z | Y)\}$$

· D-separation is sound and complete

13

Example:





• Complete the I(G) of this graph:

Eric Xino

Towards quantitative specification of probability distribution



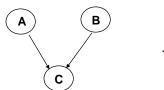
- Separation properties in the graph imply independence properties about the associated variables
- For the graph to be useful, any conditional independence properties we can derive from the graph should hold for the probability distribution that the graph represents
- The Equivalence Theorem

For a graph G,

Let \mathcal{D}_1 denote the family of all distributions that satisfy I(G), Let \mathcal{D}_2 denote the family of all distributions that factor according to G, Then $\mathcal{D}_1 \equiv \mathcal{D}_2$.

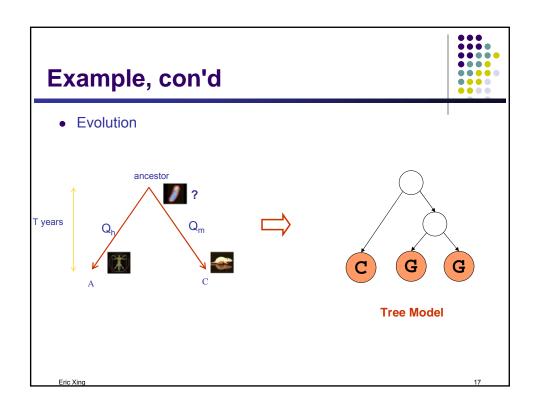
Example

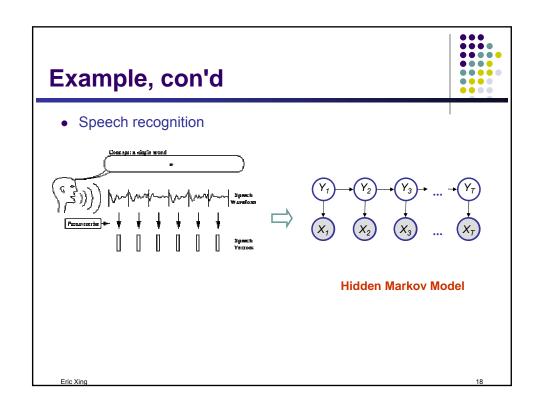


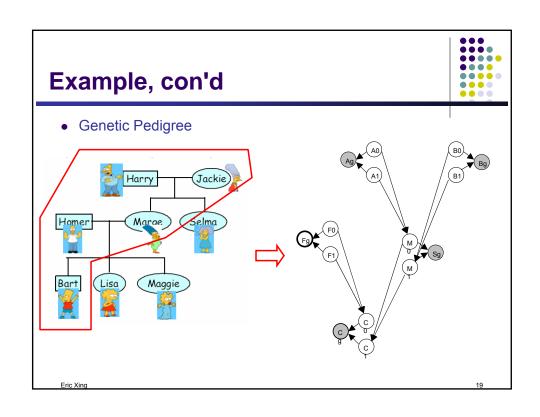


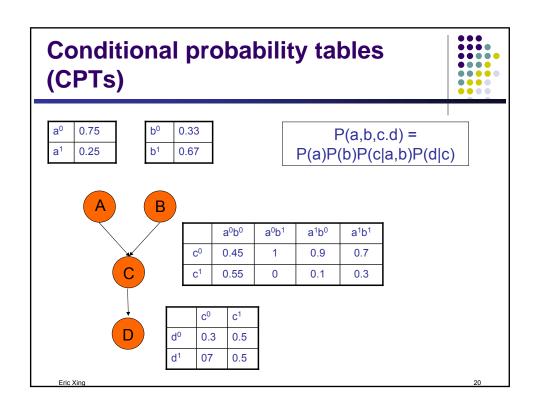
p(A,B,C) =

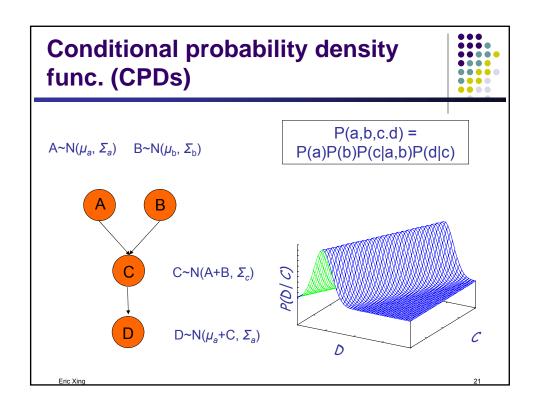
Eric Xing

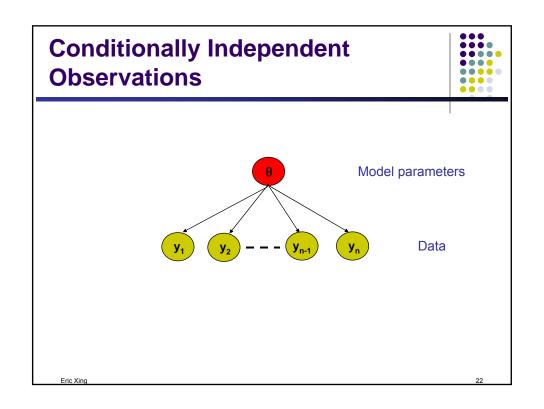






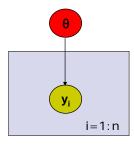






"Plate" Notation





Model parameters

Data =
$$\{y_1, ..., y_n\}$$

Plate = rectangle in graphical model

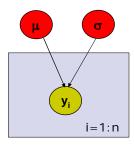
variables within a plate are replicated in a conditionally independent manner

Eric Xing

23

Example: Gaussian Model





Generative model:

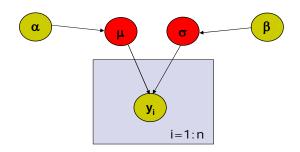
$$\begin{split} p(y_1, \dots y_n \mid \mu, \, \sigma) &= P \ p(y_i \mid \mu, \, \sigma) \\ &= \ p(\text{data} \mid \text{parameters}) \\ &= \ p(D \mid \theta) \\ \text{where } \theta &= \{\mu, \, \sigma\} \end{split}$$

- Likelihood = p(data | parameters)= p(D | θ)= L (θ)
- Likelihood tells us how likely the observed data are conditioned on a particular setting of the parameters
 - Often easier to work with log L (θ)

Eric Xing

Example: Bayesian Gaussian Model





Note: priors and parameters are assumed independent here

Eric Xin

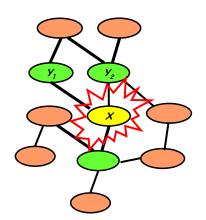
25

Markov Random Fields



Structure: an *undirected graph*

- Meaning: a node is conditionally independent of every other node in the network given its Directed neighbors
- Local contingency functions (potentials) and the cliques in the graph completely determine the joint dist.
- Give correlations between variables, but no explicit way to generate samples

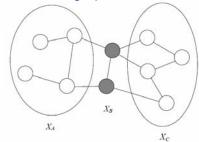


Eric Xin

Semantics of Undirected Graphs



• Let *H* be an undirected graph:



- B separates A and C if every path from a node in A to a node in C passes through a node in B: sep_H(A; C|B)
- A probability distribution satisfies the *global Markov property* if for any disjoint A, B, C, such that B separates A and C, A is independent of C given B: $I(H) = \{A \perp C | B\} : sep_H(A; C | B)\}$

Fric Xina

Representation



Defn: an undirected graphical model represents a distribution P(X₁,...,X_n) defined by an undirected graph H, and a set of positive potential functions y_c associated with cliques of H, s.t.

$$P(x_1, ..., x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(\mathbf{x}_c)$$

where *Z* is known as the partition function:

$$Z = \sum_{x_1, \dots, x_n} \prod_{c \in C} \psi_c(\mathbf{x}_c)$$

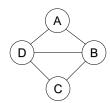
- Also known as Markov Random Fields, Markov networks ...
- The potential function can be understood as an contingency function of its arguments assigning "pre-probabilistic" score of their joint configuration.

Eric Xing

Cliques



- For $G=\{V,E\}$, a complete subgraph (clique) is a subgraph $G'=\{V'|V,E'|E\}$ such that nodes in V' are fully interconnected
- A (maximal) clique is a complete subgraph s.t. any superset V"ÉV' is not complete.
- A sub-clique is a not-necessarily-maximal clique.



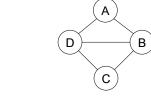
- Example:
 - max-cliques = {A,B,D}, {B,C,D},
 - sub-cliques = $\{A,B\}$, $\{C,D\}$, ... \rightarrow all edges and singletons

Eric Xing

20

Example UGM – using max cliques





$$P(x_1, x_2, x_3, x_4) = \frac{1}{Z} \psi_c(\mathbf{x}_{124}) \times \psi_c(\mathbf{x}_{234})$$



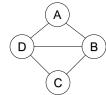
$$Z = \sum_{x_1, x_2, x_3, x_4} \psi_c(\mathbf{x}_{124}) \times \psi_c(\mathbf{x}_{234})$$

• For discrete nodes, we can represent $P(X_{1:4})$ as two 3D tables instead of one 4D table

Eric Xin

Example UGM – using subcliques





$$P(x_1, x_2, x_3, x_4) = \frac{1}{Z} \prod_{ij} \psi_{ij}(\mathbf{x}_{ij})$$



$$=\frac{1}{Z}\psi_{12}(\mathbf{x}_{12})\psi_{14}(\mathbf{x}_{14})\psi_{23}(\mathbf{x}_{23})\psi_{24}(\mathbf{x}_{24})\psi_{34}(\mathbf{x}_{34})$$

$$Z = \sum_{x_1, x_2, x_3, x_4} \prod_{ij} \psi_{ij}(\mathbf{x}_{ij})$$

 For discrete nodes, we can represent P(X_{1:4}) as 5 2D tables instead of one 4D table

Eric Xing

31

Interpretation of Clique Potentials





• The model implies $X \perp Z \mid Y$. This independence statement implies (by definition) that the joint must factorize as:

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

- We can write this as: p(x,y,z) = p(x,y)p(z|y)p(x,y,z) = p(x|y)p(z,y), but
 - cannot have all potentials be marginals
 - cannot have all potentials be conditionals
- The positive clique potentials can only be thought of as general "compatibility", "goodness" or "happiness" functions over their variables, but not as probability distributions.

Eric Xing

Exponential Form



• Constraining clique potentials to be positive could be inconvenient (e.g., the interactions between a pair of atoms can be either attractive or repulsive). We represent a clique potential $\psi_{c}(\mathbf{x}_{c})$ in an unconstrained form using a real-value "energy" function $\phi_{c}(\mathbf{x}_{c})$:

$$\psi_c(\mathbf{x}_c) = \exp\{-\phi_c(\mathbf{x}_c)\}$$

For convenience, we will call $\phi_c(\mathbf{x}_c)$ a potential when no confusion arises from the context.

This gives the joint a nice additive strcuture

$$p(\mathbf{x}) = \frac{1}{Z} \exp \left\{ -\sum_{c \in C} \phi_c(\mathbf{x}_c) \right\} = \frac{1}{Z} \exp \left\{ -H(\mathbf{x}) \right\}$$

where the sum in the exponent is called the "free energy":

$$H(\mathbf{x}) = \sum \phi_c(\mathbf{x}_c)$$

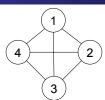
- In physics, this is called the "Boltzmann distribution".
- In statistics, this is called a log-linear model.

Eric Xing

33

Example: Boltzmann machines





• A fully connected graph with pairwise (edge) potentials on binary-valued nodes (for $x_i \in \{-1,+1\}$ or $x_i \in \{0,1\}$) is called a Boltzmann machine

$$P(x_1, x_2, x_3, x_4) = \frac{1}{Z} \exp \left\{ \sum_{ij} \phi_{ij}(x_i, x_j) \right\}$$
$$= \frac{1}{Z} \exp \left\{ \sum_{ij} \theta_{ij} x_i x_j + \sum_{i} \alpha_i x_i + C \right\}$$

• Hence the overall energy function has the form:

$$H(x) = \sum_{ij} (x_i - \mu)\Theta_{ij}(x_j - \mu) = (x - \mu)^T \Theta(x - \mu)$$

Eric Xin

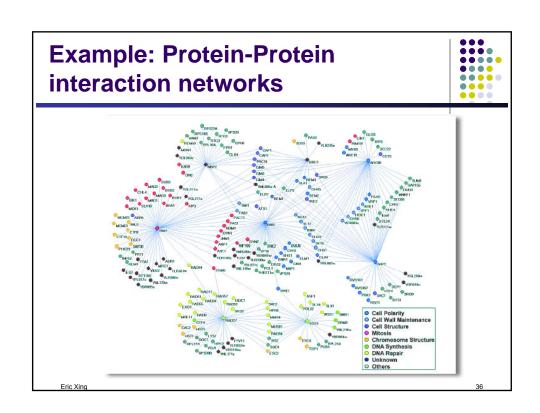
Example: Ising (spin-glass) models

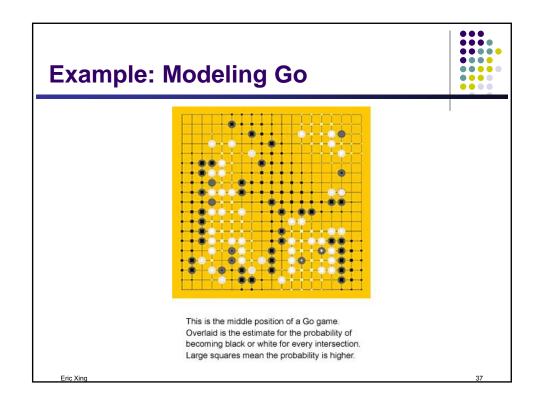


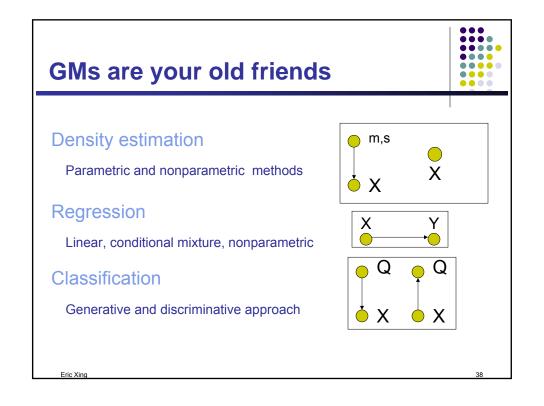
 Nodes are arranged in a regular topology (often a regular packing grid) and connected only to their geometric neighbors.

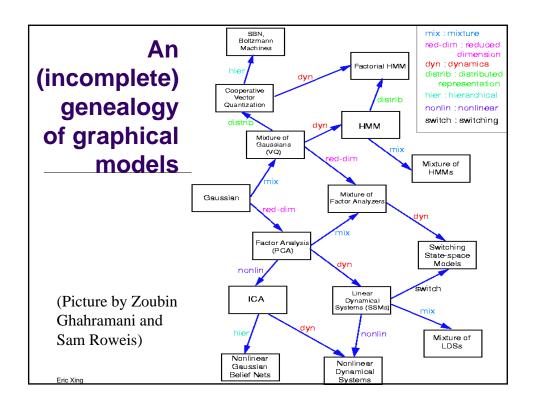
- Same as sparse Boltzmann machine, where $\theta_{ij} \neq 0$ iff i,j are neighbors.
 - e.g., nodes are pixels, potential function encourages nearby pixels to have similar intensities.
- Potts model: multi-state Ising model.

Eric Xing







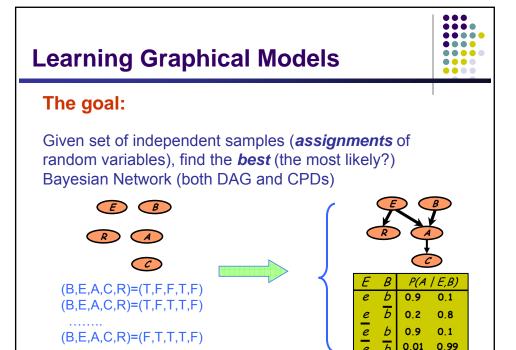


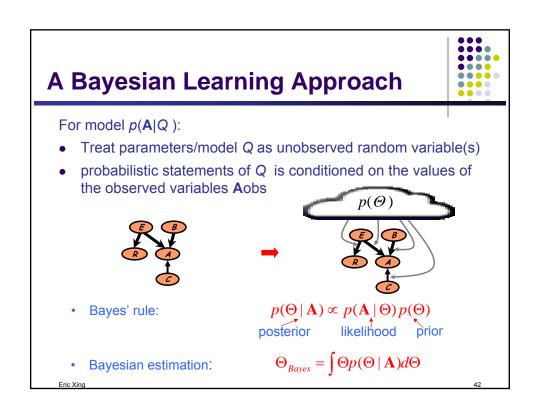
Probabilistic Inference



- Computing statistical queries regarding the network, e.g.:
- Is node X independent on node Y given nodes Z,W?
 - What is the probability of X=true if (Y=false and Z=true)?
 - What is the joint distribution of (X,Y) if Z=false?
 - What is the likelihood of some full assignment?
 - What is the most likely assignment of values to all or a subset the nodes of the network?
- General purpose algorithms exist to fully automate such computation
 - Computational cost depends on the topology of the network
 - Exact inference:
 - The junction tree algorithm
 - · Approximate inference;
 - Loopy belief propagation, variational inference, Monte Carlo sampling

Eric Xin





Why graphical models



- Probability theory provides the glue whereby the parts are combined, ensuring that the system as a whole is consistent, and providing ways to interface models to data.
- The graph theoretic side of graphical models provides both an intuitively
 appealing interface by which humans can model highly-interacting sets of
 variables as well as a data structure that lends itself naturally to the design of
 efficient general-purpose algorithms.
- Many of the classical multivariate probabilistic systems studied in fields such as statistics, systems engineering, information theory, pattern recognition and statistical mechanics are special cases of the general graphical model formalism

examples include mixture models, factor analysis, hidden Markov models, Kalman filters and Ising models.

 The graphical model framework provides a way to view all of these systems as instances of a common underlying formalism.

--- M. Jordan

ric Xing