

Introduction to PGMs: Discrete Variables

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

1. What are graphical models (or PGMs)
2. Use of PGMs Engineering and AI
3. Directionality in graphs
4. Bayesian Networks
5. Generative Models and Sampling
6. Using PGMs with fully Bayesian Models
7. Discrete Case
8. Complexity Issues

Discrete Variables

- When constructing more complex probability distributions from simpler (exponential) distributions graphical models are useful
- Graphical models have nice properties when each parent-child pair are conjugate
- Two cases of interest:
 - Both correspond to discrete variables
 - Both correspond to Gaussian variables

Discrete Case

- Probability distribution $p(\mathbf{x} | \boldsymbol{\mu})$ for a single discrete variable x having K states



- Using 1-of- K representation

– For $K=6$ when $x_3=1$ then x represented as $\mathbf{x}=(0,0,1,0,0,0)^T$

Note that $\sum_{k=1}^K x_k = 1$

- If probability of $x_k=1$ is given by parameter μ_k

$$p(\mathbf{x} | \boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad \text{where } \boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T$$

- The distribution is normalized: $\sum_{k=1}^K \mu_k = 1$

– There are $K-1$ independent values for μ_k needed to define distribution

Two Discrete Variables

$x_1:$	x_{11}					x_{1K}
$x_2:$	x_{21}					x_{2K}

- x_1 and x_2 each with K states each
- Denote probability of both $x_{1k}=1$ and $x_{2l}=1$ by μ_{kl}
 - where x_{1k} denotes k^{th} component of x_1

- Joint distribution is

$$p(\mathbf{x}_1, \mathbf{x}_2 | \boldsymbol{\mu}) = \prod_{k=1}^K \prod_{l=1}^K \mu_{kl}^{x_{1k} x_{2l}}$$

- Since parameters are subject to constraint $\sum_k \sum_l \mu_{kl} = 1$
- There are $K^2 - 1$ parameters
- For arbitrary distribution over M variables there are $K^M - 1$ parameters

Graphical Models for Two Discrete Variables

- Joint distribution $p(\mathbf{x}_1, \mathbf{x}_2)$
- Using product rule
 - is factored as $p(\mathbf{x}_2|\mathbf{x}_1)p(\mathbf{x}_1)$
- Has two node graph
- Marginal distribution $p(\mathbf{x}_1)$ has $K-1$ parameters
- Conditional distribution $p(\mathbf{x}_2|\mathbf{x}_1)$ also requires $K-1$ parameters for each of K values of \mathbf{x}_1
- Total number of parameters is $(K-1) + K(K-1) = K^2 - 1$
 - As before



Two Independent Discrete Variables

- x_1 and x_2 are independent
 - Has graphical model
- Each variable described by a separate multinomial distribution
 - Total no of parameters is $2(K-1)$
 - For M variables no of parameters is $M(K-1)$
- Reduced number of parameters by dropping links in graph
 - Grows linearly with no of variables



Fully connected has high complexity

- General case of M discrete variables X_1, \dots, X_M
- If BN is fully connected
 - Completely general distribution with K^{M-1} parameters
- If there are no links
 - Joint distribution factorizes into product of marginals
 - Total no of parameters is $M(K-1)$
- Graphs of intermediate levels of connectivity
 - More general distribution than fully factorized ones
 - Require fewer parameters than general joint distribution
 - Example: chain of nodes

Special Case: Chain of Nodes



- Marginal distribution $p(x_1)$ requires $K-1$ parameters
- Each of the $M-1$ conditional distributions $p(x_i|x_{i-1})$, for $i=2,\dots,M$ requires $K(K-1)$ parameters
- Total parameter count is
$$K-1+(M-1)K(K-1)$$

Which is quadratic in K
- Grows linearly (not exponentially) with length of chain

Alternative: Sharing Parameters

- Reduce parameters by *sharing* or *tying* parameters

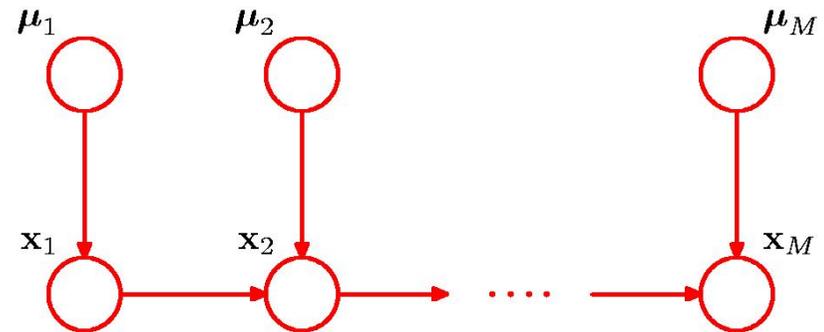


- In above,
 - all conditional distributions $p(x_i|x_{i-1})$, for $i=2,\dots,M$
 - share same set of $K(K-1)$ parameters governing distribution of x_1
- Total of K^2-1 parameters needed to specify distribution

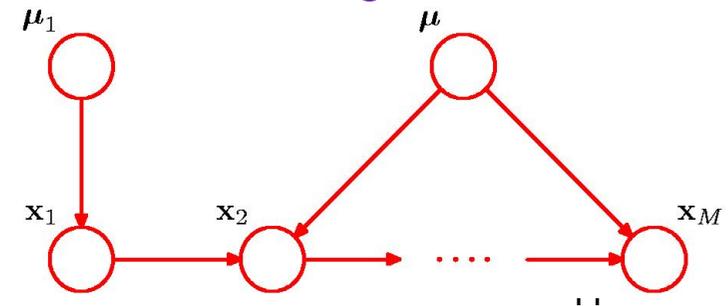
Conversion into Bayesian Model

- Given graph over discrete variables
- We can turn it into a Bayesian model by introducing Dirichlet priors for parameters
- Each node acquires an additional parent for each discrete node
- Tie the parameters governing conditional distributions $p(\mathbf{x}_i | \mathbf{x}_{i-1})$

Chain of Nodes
With priors



Sharing Parameters



Binomial: Beta Prior

- Bernoulli:

$$p(x=1|\mu) = \mu$$

- Likelihood of Bernoulli with $D = \{x_1, \dots, x_N\}$

$$p(D|\mu) = \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n}$$

- Binomial:

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}$$

- Conjugate Prior:

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

Multinomial: Dirichlet Prior

- **Generalized Bernoulli (1-of- K)** $\mathbf{x}=(0,0,1,0,0,0)^T$ $K=6$

$$p(\mathbf{x} | \boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \text{ where } \boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T \quad K=2 \text{ is Bernoulli}$$

- **Multinomial** $K=2$ is Binomial

$$Mult(m_1 m_2 \dots m_K | \boldsymbol{\mu}, N) = \binom{N}{m_1 m_2 \dots m_K} \prod_{k=1}^K \mu_k^{m_k} \quad \binom{N}{m_1 m_2 \dots m_k} = \frac{N!}{m_1! m_2! \dots m_k!}$$

- Where the normalization coefficient is the no of ways of partitioning N objects into K groups of size $m_1, m_2 \dots m_k$
- **Conjugate prior distribution for parameters m_k**

$$p(\boldsymbol{\mu} | \boldsymbol{\alpha}) \propto \prod_{k=1}^K \mu_k^{\alpha_k - 1} \text{ where } 0 \leq \mu_k \leq 1 \text{ and } \sum_k \mu_k = 1$$

- **Normalized form is**

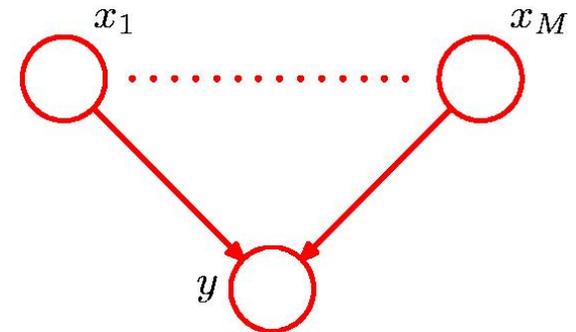
$$Dir(\boldsymbol{\mu} | \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} \prod_{k=1}^K \mu_k^{\alpha_k - 1} \text{ where } \alpha_0 = \sum_{k=1}^K \alpha_k$$

Controlling Number of parameters in models: Parameterized Conditional Distributions

- Control exponential growth of parameters in models of discrete variables
- Use parameterized models for conditional distributions instead of complete tables of conditional probability values

Parameterized Conditional Distributions

- Consider graph with binary variables
- Each parent variable x_i governed by single parameter μ_i representing probability $p(x_i=1)$
- M parameters in total for parent nodes
- Conditional distribution $p(y|x_1, \dots, x_M)$ requires 2^M parameters
 - Representing probability $p(y=1)$ for each of the 2^M settings of parent variables
 - 000000000 to 111111111



Conditional distribution using logistic sigmoid

- Parsimonious form of conditional distribution
- Logistic sigmoid acting on linear combination of parent variables

$$p(y = 1 | x_1, \dots, x_M) = \sigma \left(w_0 + \sum_{i=1}^M w_i x_i \right) = \sigma(\mathbf{w}^T \mathbf{x})$$

– **where** $\sigma(a) = (1 + \exp(-a))^{-1}$ is the logistic sigmoid

– $\mathbf{x} = (x_0, x_1, \dots, x_M)^T$ is vector of parent states

- No of parameters grows linearly with M
- Analogous to choice of a restrictive form of covariance matrix in multivariate Gaussian

Linear Gaussian Models

- Expressing multivariate Gaussian as a directed graph
- Corresponding to linear Gaussian model over component variables
 - Mean of a conditional distribution is a linear function of the conditioning variable
- Allows expressing interesting structure of distribution
 - General Gaussian case and diagonal covariance case represent opposite extremes

Graph with continuous random variables

- Arbitrary acyclic graph over D variables
- Node i represents a single continuous random variable x_i having Gaussian distribution
- Mean of distribution is a linear combination of states of its parent nodes pa_i of node i

$$p(x_i | pa_i) = \mathcal{N}\left(x_i \mid \sum_{j \in pa_i} w_{ij} x_j + b_i, v_i\right)$$

- Where
 - w_{ij} and b_i are parameters governing the mean
 - v_i is the variance of the conditional distribution for x_i

Joint Distribution

- Log of joint distribution

$$\begin{aligned}\ln p(\mathbf{x}) &= \sum_{i=1}^D \ln p(x_i | pa_i) \\ &= -\sum_{i=1}^D \frac{1}{2v_i} \left(x_i - \sum_{j \in pa_i} w_{ij} x_j - b_i \right)^2 + \underbrace{const}_{\text{Terms independent of } \mathbf{x}}\end{aligned}$$

Terms independent of \mathbf{x}

- Where $\mathbf{x} = (x_1, \dots, x_D)^T$
- This is a quadratic function of \mathbf{x}
 - Hence joint distribution $p(\mathbf{x})$ is a multivariate Gaussian

Mean and Covariance of Joint Distribution

- Recursive Formulation

- Since each variable x_i has,

- conditional on the states of its parents,

- a Gaussian distribution, we can write $x_i = \sum_{j \in pa_i} w_{ij} + b_i + \sqrt{v_i} \varepsilon_i$

- where ε_i is

- a zero mean, unit variance Gaussian random variable
- satisfying $\mathbf{E}[\varepsilon_i] = 0$ and $\mathbf{E}[\varepsilon_i \varepsilon_j] = I_{ij}$
- and I_{ij} is the i, j element of the identity matrix

- Taking expectation $\mathbf{E}[x_i] = \sum_{j \in pa_i} w_{ij} \mathbf{E}[x_j] + b_i$

- Thus we can find components of $\mathbf{E}[\mathbf{x}] = (\mathbf{E}[x_1], \dots, \mathbf{E}[x_D])^T$

- by starting at lowest numbered node and working recursively through the graph

- Similarly elements of covariance matrix

$$\text{cov}[x_i, x_j] = \sum_{k \in pa_j} w_{jk} \text{cov}[x_i, x_k] + I_{ij} v_j$$

Three cases for no. of parameters

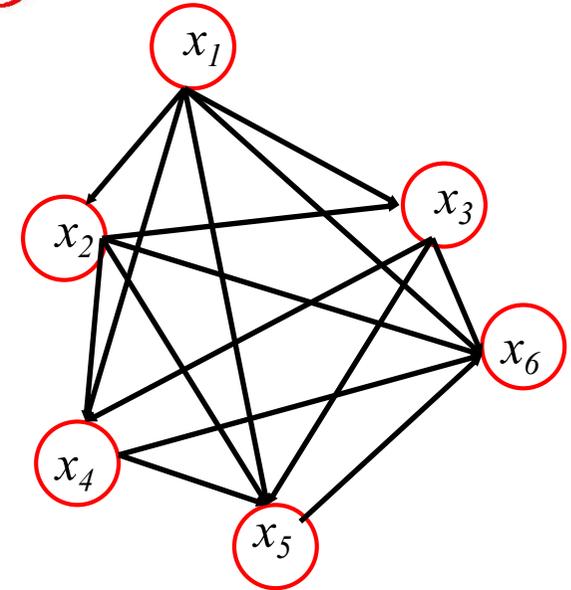
- No links in the graph

– $2D$ parameters



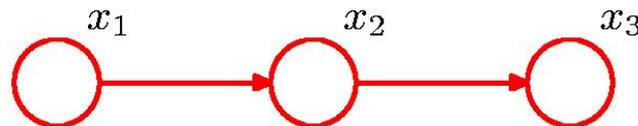
- Fully connected graph

– $D(D+1)/2$ parameters



- Graphs with intermediate level of complexity

– Chain



Extreme Case with no links

- D isolated nodes
 - There are no parameters w_{ij}
 - Only D parameters b_i and
 - D parameters v_i
- Mean of $p(x)$ given by $(b_1, \dots, b_D)^T$
- Covariance matrix is diagonal of form $\text{diag}(v_1, \dots, v_D)$
- Joint distribution has total of $2D$ parameters
- Represents set of D independent univariate Gaussian distributions

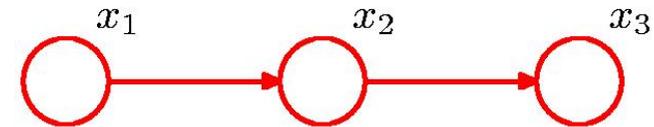
Extreme case with all links

- Fully connected graph
- Each node has all lower numbered nodes as parents
- Matrix w_{ij} has $i-1$ entries on the i th row and hence is a lower triangular matrix (with no entries on leading diagonal)
- Total no of parameters w_{ij} is to take D^2 no of elements in $D \times D$ matrix, subtracting D to account for diagonal and divide by 2 to account for elements only below diagonal

Graph with intermediate complexity

- Link missing between variables x_1 and x_3
- Mean and covariance of joint distribution are

$$\mu = (b_1, b_2 + w_{21}b_1, b_3 + w_{32}b_2 + w_{32}w_{21}b_1)^T$$



$$\Sigma = \begin{pmatrix} v_1 & w_{21}v_1 & w_{32}w_{21}v_1 \\ w_{21}v_1 & v_2 + w_{21}^2v_1 & w_{32}(v_2 + w_{21}^2v_1) \\ w_{32}w_{21}v_1 & w_{32}(v_2 + w_{21}^2v_1) & w_{32}^2(v_2 + w_{21}^2v_1) \end{pmatrix}$$

Extension to multivariate Gaussian variables

- Nodes in the graph represent multivariate Gaussian variables
- Write conditional distribution for node i in the form

$$p(\mathbf{x}_i | pa_i) = \mathcal{N}\left(\mathbf{x}_i \left| \sum_{j \in pa_i} \mathbf{W}_{ij} \mathbf{x}_j + \mathbf{b}_i, \sum v_i \right.\right)$$

- where \mathbf{W}_{ij} is a matrix (non-square if \mathbf{x}_i and \mathbf{x}_j have different dimensionalities)

Summary

1. PGMs allow visualizing probabilistic models
 - Joint distributions are directed/undirected PGMs
2. PGMs can be used to generate samples
 - Ancestral sampling with directed PGMs is simple
3. PGMs are useful for Bayesian statistics
 - Discrete variable PGM represented using Dirichlet priors
4. Parameter explosion controlled by tying parameters
5. Multivariate Gaussian expressed as PGM
 - Graph is a linear Gaussian model over components