

# Gaussian Bayesian Networks

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

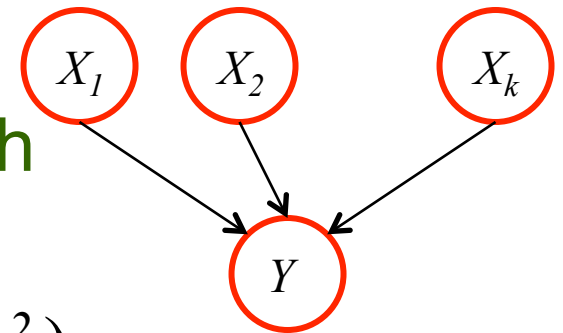
- Linear Gaussian Model
- Learning Parameters of a Gaussian Bayesian Network
- Multivariate Gaussian representations
  1. Covariance Form
  2. Information Form:
    1. Bayesian Networks
    2. Markov Random Fields

# Linear Gaussian Model

- $Y$  has a linear Gaussian model with continuous parents  $X_1, \dots, X_k$

with parameters  $\beta_0, \dots, \beta_k$  and  $\sigma^2$  such that

$$P(Y|X_1, \dots, X_k) \sim N(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k; \sigma^2)$$



- In vector notation

$$P(Y|\mathbf{X}) \sim N(\beta_0 + \beta^T \mathbf{X}; \sigma^2)$$

- Can be viewed as  $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon$ 
  - i.e.,  $Y$  is a linear function of the variables with the addition of Gaussian noise with mean 0 and variance  $\sigma^2$

Linear regression could be used to solve for the parameters

# Gaussian Bayesian Network

- *Definition:* All variables are Gaussian and all CPDs are linear Gaussian
- Let  $Y$  be a linear Gaussian of its parents  $X_1, \dots, X_k$ 
  - i.e., mean of  $Y$  is a linear combination of means of Gaussian parents

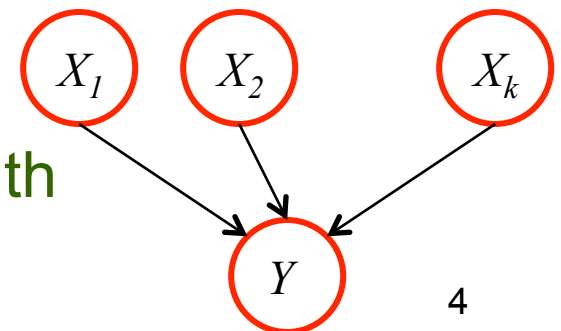
$$P(Y|X) \sim N(\beta_0 + \beta^t X; \sigma^2)$$

- If  $X_1, \dots, X_k$  are jointly normal  $N(\mu; \Sigma)$  then the distribution of  $Y$  is normal  $N(\mu_Y; \sigma_Y^2)$  where

$$\mu_Y = \beta_0 + \beta^t \mu \text{ and } \sigma_Y^2 = \sigma^2 + \beta^t \Sigma \beta$$

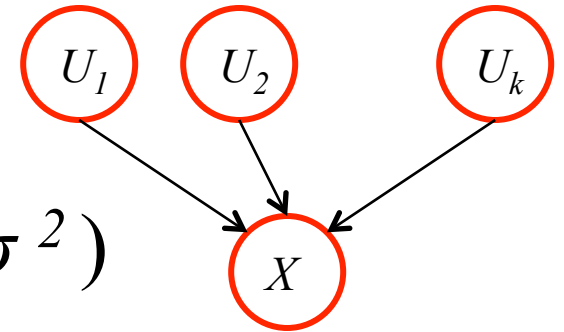
The joint distribution over  $(X, Y)$  is normal with

$$\text{Cov}[X_i : Y] = \sum_{j=1}^k \beta_j \Sigma_{i,j}$$



# Learning parameters of Gaussian Bayesian Network

- All variables are Gaussian
- $P(X|u_1, \dots, u_k) \sim N(\beta_0 + \beta_1 u_1 + \dots + \beta_k u_k; \sigma^2)$
- Task is to learn the parameters



$$\theta_{X|U} = \{\beta_0, \beta_1, \dots, \beta_k, \sigma^2\}$$

- To find ML estimate of parameters, define log-likelihood

$$\log L_X(\theta_{X|U} : D) = \sum_m \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\beta_0 + \beta_1 u_1[m] + \dots + \beta_k u_k[m] - x[m])^2 \right]$$

- Closed-form solutions for parameters

# Solving for the parameters

- Gradient of log-likelihood wrt  $\beta_0$

$$-\frac{1}{\sigma^2} \left( M\beta_0 + \beta_1 \sum_m u_1[m] + \dots + \beta_k \sum_m u_k[m] - x[m] \right)$$

- Equating to zero and rearranging we get

$$\frac{1}{M} \sum_m x[m] = \beta_0 + \beta_1 \frac{1}{M} \sum_m u_1[m] + \dots + \beta_k \frac{1}{M} \sum_m u_k[m]$$

- All the summations can be obtained from data thus giving us a linear equation
- Similarly we get  $k$  more linear equations by taking derivatives wrt  $\beta_i$
- Standard linear algebra techniques are used to solve  $k+1$  simultaneous equations <sup>6</sup>

# Solving Linear Equations

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2 \\ \vdots & \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n &= b_m \end{aligned}$$

$$x_1 \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{m1} \end{bmatrix} + x_2 \begin{bmatrix} a_{12} \\ a_{22} \\ \vdots \\ a_{m2} \end{bmatrix} + \cdots + x_n \begin{bmatrix} a_{1n} \\ a_{2n} \\ \vdots \\ a_{mn} \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

1. Matrix Solution:  $\mathbf{Ax}=\mathbf{b}$   
Therefore  $\mathbf{x}=\mathbf{A}^{-1}\mathbf{b}$

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

2. Gaussian Elimination

$$\begin{cases} x + 3y - 2z = 5 \\ 3x + 5y + 6z = 7 \\ 2x + 4y + 3z = 8 \end{cases}$$

$$\begin{aligned} \left[ \begin{array}{ccc|c} 1 & 3 & -2 & 5 \\ 3 & 5 & 6 & 7 \\ 2 & 4 & 3 & 8 \end{array} \right] &\sim \left[ \begin{array}{ccc|c} 1 & 3 & -2 & 5 \\ 0 & -4 & 12 & -8 \\ 2 & 4 & 3 & 8 \end{array} \right] &\sim \left[ \begin{array}{ccc|c} 1 & 3 & -2 & 5 \\ 0 & -4 & 12 & -8 \\ 0 & -2 & 7 & -2 \end{array} \right] &\sim \left[ \begin{array}{ccc|c} 1 & 3 & -2 & 5 \\ 0 & 1 & -3 & 2 \\ 0 & -2 & 7 & -2 \end{array} \right] \\ &\sim \left[ \begin{array}{ccc|c} 1 & 3 & -2 & 5 \\ 0 & 1 & -3 & 2 \\ 0 & 0 & 1 & 2 \end{array} \right] &\sim \left[ \begin{array}{ccc|c} 1 & 3 & -2 & 5 \\ 0 & 1 & 0 & 8 \\ 0 & 0 & 1 & 2 \end{array} \right] &\sim \left[ \begin{array}{ccc|c} 1 & 3 & 0 & 9 \\ 0 & 1 & 0 & 8 \\ 0 & 0 & 1 & 2 \end{array} \right] &\sim \left[ \begin{array}{ccc|c} 1 & 0 & 0 & -15 \\ 0 & 1 & 0 & 8 \\ 0 & 0 & 1 & 2 \end{array} \right] \end{aligned}$$

# Estimating variance

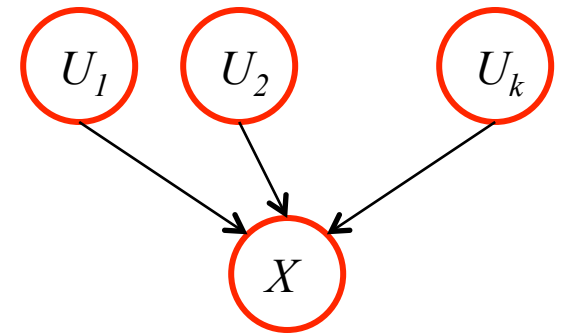
- Taking derivative of likelihood and setting to zero, we get

$$\sigma^2 = Cov_D[X; X] - \sum_i \sum_j \beta_i \beta_j Cov_D[U_i; U_j]$$

– where

$$Cov_D[X; Y] = E_D[X \cdot Y] - E_D[X] \cdot E_D[Y]$$

$$E_D[X \cdot Y] = \frac{1}{M} \sum_m x[m]y[m] \quad E_D[X] = \frac{1}{M} \sum_m x[m]$$



- First term is the empirical variance of  $X$
- Other terms are empirical covariances of inputs

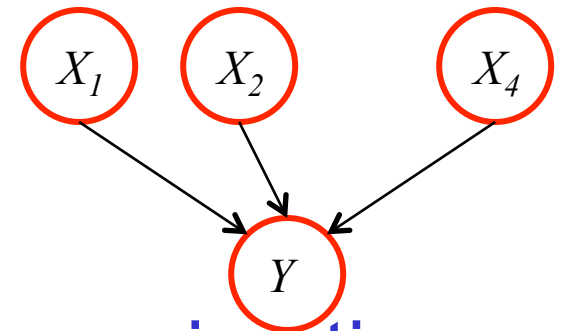


# Properties of Linear Gaussian

- *Advantage*: Simple model that captures many dependencies
- *Disadvantage*: cannot capture dependence of variance of child on values of parents

– Can be extended

- E.g., mean of  $Y$  is  $\sin(x_1)^{x_2}$   
variance is  $(x_3/x_4)^2$



- But linear Gaussian is a good approximation
- Provides an alternative for representing multivariate Gaussian distributions

# Multivariate Gaussian

- Three equivalent representations

1. Basic parameterization (covariance form)

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

$$\begin{aligned} \boldsymbol{\mu} &= E[\mathbf{x}] \\ \Sigma &= E[\mathbf{x}\mathbf{x}^t] - E[\mathbf{x}]E[\mathbf{x}]^t \end{aligned}$$

$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ -3 \\ 4 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ -2 & -5 & 8 \end{bmatrix}$$

$x_1$  is negatively correlated with  $x_3$   
 $x_2$  is negatively correlated with  $x_3$

2. Gaussian Bayesian network (information form)

- Conversion to Information form (precision matrix) captures conditional independencies

3. Gaussian Markov Random Field

- Easiest conversion

# Gaussian Distribution Operations

- There are two main operations we wish to perform on a distribution:
  1. Marginal distribution over some subset  $Y$
  2. Conditioning the distribution on some  $Z=z$
- Marginalization is trivially performed in the covariance form
- Conditioning a Gaussian is very easy to perform in the information form

# Marginalization in Covariance Form

- Can be simply read from the mean and covariance matrix
- Decompose the mean and covariance as:

$$p(\mathbf{X}, \mathbf{Y}) = N\left(\begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}; \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix}\right)$$

- Marginal distribution over  $\mathbf{Y}$  is a normal distribution  $N(\mu_Y; \Sigma_{YY})$

# Determining Conditional Distribution using Covariance Form

- Let  $(\mathbf{X}, Y)$  have a normal distribution defined by

$$p(\mathbf{X}, Y) = N\left(\begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}; \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix}\right)$$

- Then marginal distribution  $P(Y|\mathbf{X})$  is given by

$P(Y|\mathbf{X}) \sim N(\beta_0 + \beta^T \mathbf{X}; \sigma^2)$  where

$$\begin{aligned} \beta_0 &= \mu_Y - \Sigma_{YX} \Sigma_{XX}^{-1} \mu_X \\ \beta &= \Sigma_{XX}^{-1} \Sigma_{YX} \\ \sigma^2 &= \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY} \end{aligned}$$

- This allows us to take a joint Gaussian distribution and construct a Bayesian network

# Constructing a BN from a Gaussian

- Definition of I-map
  - $I$  is set of independence assertions of form  $(X \perp Y / Z)$ , which are true of distribution  $P$
  - $\mathcal{G}$  is an I-map of  $P$  implies:  $I_{\mathcal{G}} \subseteq I$ 
    - Minimal I-map: removal of a single edge renders the graph not an I-map
- Algorithm uses the principle that for  $Z$  to be parents of node  $X$ ,  $(X \perp Y / Z)$  where  $Y = \chi - Z$ 
  - need to check if  $P(X/Z)$  is independent of  $P(Y_i/Z)$ 
    - Can be done using deviance measure:  $\chi^2$  or mutual info (KL-div)
 

$$d_{\chi^2}(\mathcal{D}) = \sum_{x_i, x_j} \frac{(M[x_i, x_j] - M \cdot \hat{P}(x_i) \cdot \hat{P}(x_j))^2}{M \cdot \hat{P}(x_i) \cdot \hat{P}(x_j)}$$

$$d_I(\mathcal{D}) = \frac{1}{M} \sum_{x_i, x_j} M[x_i, x_j] \log \frac{M[x_i, x_j]}{M[x_i]M[x_j]}$$
    - Simpler Method: use Information form of Gaussian

# BN from conditional independencies

Algorithm computes an I-Map which may not be unique

- $\mathcal{X} = \{X_1, \dots, X_n\}$ : set of random variables
- $\mathcal{I}$ : set of independencies

```

1  Set  $\mathcal{G}$  to an empty graph over  $\mathcal{X}$ 
2  for  $i = 1, \dots, n$ 
3     $U \leftarrow \{X_1, \dots, X_{i-1}\}$  //  $U$  is the current candidate for parents of  $X_i$ 
4    for  $U' \subseteq \{X_1, \dots, X_{i-1}\}$ 
5      if  $U' \subset U$  and  $(X_i \perp \{X_1, \dots, X_{i-1}\} - U' \mid U') \in \mathcal{I}$  then
6         $U \leftarrow U'$ 
7      // At this stage  $U$  is a minimal set satisfying  $(X_i \perp$ 
8         $\{X_1, \dots, X_{i-1}\} - U \mid U)$ 
9      // Now set  $U$  to be the parents of  $X_i$ 
10     for  $X_j \in U$ 
11       Add  $X_j \rightarrow X_i$  to  $\mathcal{G}$ 
12  return  $\mathcal{G}$ 

```

Nodes other than  $U'$

# Information form gives edges of BN

- Standard Form

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

- Using Precision  $J = \Sigma^{-1}$

$$-\frac{1}{2}(x - \mu)^t \Sigma^{-1}(x - \mu) = -\frac{1}{2}(x - \mu)^t J(x - \mu) = -\frac{1}{2}\left[x^t Jx - 2x^t J\mu + \mu^t J\mu\right]$$

– Since last term is constant, we get information form

$$p(x) \propto \exp\left(-\frac{1}{2}x^t Jx + (J\mu)^t x\right)$$

- $J\mu$  is called the potential vector
- $J_{i,j} = 0$  iff  $X_i$  is independent of  $X_j$  given  $\mathcal{X} - \{X_i, X_j\}$

– Example

$$\mu = \begin{bmatrix} 1 \\ -3 \\ 4 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ -2 & -5 & 8 \end{bmatrix} \quad J = \begin{bmatrix} 0.3125 & -0.125 & 0 \\ -0.125 & 0.5833 & 0.3333 \\ 0 & 0.3333 & 0.3333 \end{bmatrix}$$

- Since  $J_{1,3} = 0$ ,  $x_1$  and  $x_3$  are independent given  $x_2$

– Non-zero entries define edges (non-correlations)



# Gaussian Markov Random Fields

- Follows directly from information form

$$p(\mathbf{x}) \propto \exp\left[-\frac{1}{2}\mathbf{x}^t J \mathbf{x} + (J\boldsymbol{\mu})^t \mathbf{x}\right]$$

– which is obtained from covariance form with  $J = \Sigma^{-1}$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$$

- Break-up exponent into two types of terms

– Using the potential vector  $\mathbf{h} = J\boldsymbol{\mu}$

– Terms involving single variable  $X_i$

$$-\frac{1}{2} J_{i,i} x_i^2 + h_i x_i$$

- Called node potentials

Terms involving pairs of variables  $X_i, X_j$

$$-\frac{1}{2} \left[ J_{i,j} x_i x_j + J_{j,i} x_j x_i \right] = -J_{i,j} x_i x_j$$

- Called edge potentials (when  $J_{i,j} = 0$ , there is no edge)