

Discrete Variables & Linear Gaussian Models

Graphical models are a good way of linking together simple (exp. family (Chapt. 2.4)) distributions to form more powerful and useful models of more complicated data and processes.

Two cases are particularly elegant, discrete and lin. Gaussian (i.e. which all parent-child relations are either discrete or lin. G.)

Discrete (to start)

Code a discrete r.v. \vec{X} as $\vec{x} = \underbrace{[0 \ 0 \ 1 \ 0 \ 0 \ 0]}_{K \text{ possible states}}$

A discrete dist. with parameters $\vec{\mu}$ can be written

$$p(\vec{x} | \vec{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad \text{eg. } \vec{\mu} = [.1, .1, .3, .5, 0, 0]$$

Since $\sum \mu_k = 1$ only $K-1$ parameters are required to represent $\vec{\mu}$

With two discrete r.v.'s \vec{x}_1 and \vec{x}_2 (each w/ K states) we can write the joint distribution that $x_{1k} = 1$ and $x_{2\ell} = 1$ as

$$p(\vec{x}_1, \vec{x}_2 | \mu) = \prod_{k=1}^K \prod_{\ell=1}^K \mu_{k\ell}^{x_{1k} x_{2\ell}}$$

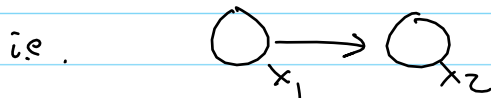
where $\sum \mu_{k\ell} = 1$ and $\mu_{k\ell}$ is a matrix with $K^2 - 1$ free parameters.

For an arbitrary joint distribution with M discrete variables, K^{M-1} params are needed

k^{M-1} grows exponentially with M , the number of variables.

Using the prod rule we can write

$$p(\vec{x}_1, \vec{x}_2) \text{ in the form } p(x_2|x_1)p(x_1)$$



$p(x_1)$ has $K-1$ params

$p(x_2|x_1)$ requires $K-1$ params for each of the K states x_1 can be in

the total # params is

$$K-1 + K(K-1) = K^2 - 1 \text{ as before}$$

but...

if $x_1 \perp\!\!\!\perp x_2$ ie. x_1 x_2

then the total number of params would be $2(K-1)$. For $M \perp\!\!\!\perp$ discrete R.V.'s with K states # params grows linearly (!)

$$M(K-1)$$

By dropping links we have decreased the number of parameters in the joint distribution.

Intermediate factorizations are interesting

Consider



The total number of params grows as

$k-1 + (M-1)(k-1)k$
 which is quadratic in k and linear in M

Parameters count can be reduced through sharing or tying of parameters. In (Chain) we can say all cond. dist's are ~~the same~~ ^{shared} leaving

$$k^2 - 1$$

parameters that must be specified in order to define the joint dist'n.

Linear Gaussian Models

D dimensional Gaussian w/ diagonal covariance has $2D$ free parameters, D for the mean and D for the diagonal of the covariance matrix.

D dimensional full Gaussian has D free parameters for the mean and $(D^2 - D)/2 + D$ free covariance matrix parameters (symmetric)

It is possible to co-struct intermediate Gaussian dist's of interm. complexity

Let

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

where w_{ij} and b_i are parameters that govern the conditional mean of x_i and v_i is the variance of x_i 's conditional dist.

By inspection of the log joint $p(\vec{X})$ we can see that the full joint dist. is Gaussian

$$\begin{aligned} \ln p(\vec{X}) &= \sum_{i=1}^D \ln p(x_i | \text{pa}_i) \\ &= \sum_{i=1}^D \frac{1}{2v_i} \left(x_i - \sum_{j \in \text{pa}_i} w_{ij} x_j - b_i \right)^2 + \text{const} \end{aligned}$$

which is quadratic in the components of \vec{X} and therefore Gaussian.

The mean and variance of the resulting ^{joint} Gaussian can be determined recursively

Each x_i , conditioned on its parents can be written as

$$x_i = \sum_{j \in \text{pa}_i} w_{ij} x_j + b_i + \sqrt{v_i} \varepsilon_i$$

where $\varepsilon_i \sim \mathcal{N}(0, 1)$, $\mathbb{E}[\varepsilon_i] = 0$ and $\mathbb{E}(\varepsilon_i \varepsilon_j) = I_{ij}$ where I_{ij} is the ij th element of the identity matrix.

$$\text{So } \mathbb{E}[x_i] = \sum_{j \in \text{pa}_i} w_{ij} \mathbb{E}[x_j] + b_i$$

which is how the mean of $\vec{X} \sim \mathcal{N}(\mu, \Sigma)$ can be calculated.

The covariance between x_i and x_j can also be recursively calculated.

Starting with the definition of covariance

$$\begin{aligned} \text{cov}[x_i, x_j] &= E[(x_i - E(x_i))(x_j - E(x_j))] \\ &= E[(x_i - E(x_i))(\sum_{k \in \text{pa}_j} w_{jk} x_k + b_j + \sqrt{v_j} \varepsilon_j \\ &\quad - (\sum_{k \in \text{pa}_j} w_{jk} E[x_k] + b_j))] \\ &= E[(x_i - E(x_i))(\sum_{k \in \text{pa}_j} w_{jk} (x_k - E[x_k]) + \sqrt{v_j} \varepsilon_j)] \end{aligned}$$

$$= \sum_{k \in \text{pa}_j} w_{jk} \text{cov}(x_i, x_k) + E[(x_i - E(x_i)) \sqrt{v_j} \varepsilon_j]$$

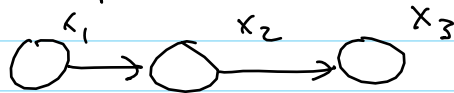
$$= E[x_i \sqrt{v_j} \varepsilon_j]$$

$$= E[(\underbrace{\sum_{j \in \text{pa}_i} w_{ij} x_j + b_i + \sqrt{v_i} \varepsilon_i}_{\text{contains no } \varepsilon \text{ terms}}) \sqrt{v_j} \varepsilon_j]$$

$$\begin{aligned} &= E[\sqrt{v_i} \sqrt{v_j} \varepsilon_i \varepsilon_j] = E[\underbrace{\varepsilon_i \varepsilon_j}_{\text{iff } i=j}] \sqrt{v_i} \sqrt{v_j} \\ &= \sqrt{v_i} \sqrt{v_j} = v_i \text{ or } v_j \end{aligned}$$

$$= \sum_{k \in \text{pa}_j} w_{jk} \text{cov}(x_i, x_k) + I_{ij} v_j$$

Example



think: latent process in Kalman filter, eg.

rules: $E[x_i] = \sum_{j \in pa_i} \omega_{ij} E[x_j] + b_i$, $cov[x_i, x_j] = \sum_{k \in pa_j} \omega_{jk} cov[x_i, x_k] + I_{ij} v_j$

$$\mu = [E[x_1], E[x_2], E[x_3]]$$

$$= [b_1, \omega_{12} b_1 + b_2, \omega_{32} (\omega_{12} b_1 + b_2) + b_3]$$

$$\sum_{i=1}^3 \begin{bmatrix} cov(x_1, x_1) & cov(x_1, x_2) & cov(x_1, x_3) \\ cov(x_2, x_1) & cov(x_2, x_2) & cov(x_2, x_3) \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$= \begin{bmatrix} v_1 & \omega_{21} v_1 & \omega_{32} v_{21} v_1 \\ & & \\ & & \end{bmatrix}$$

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