

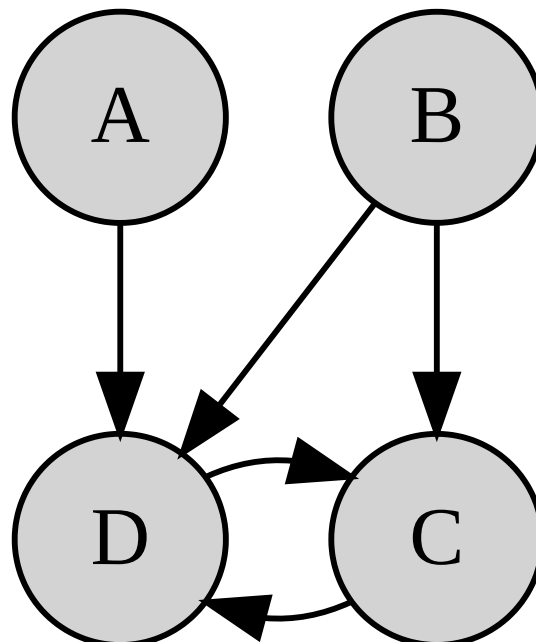
# MGM: VAE

Category	
Files	
Created	@May 2, 2023 4:24 PM
Reminder	
Status	Open
URL	
Updated	@May 5, 2023 2:02 PM

## Problem-1

(a)

Probabilistic graphical model as PGM is a joint probabilistic distribution that uses a graph structure to encode conditional independence assumptions.

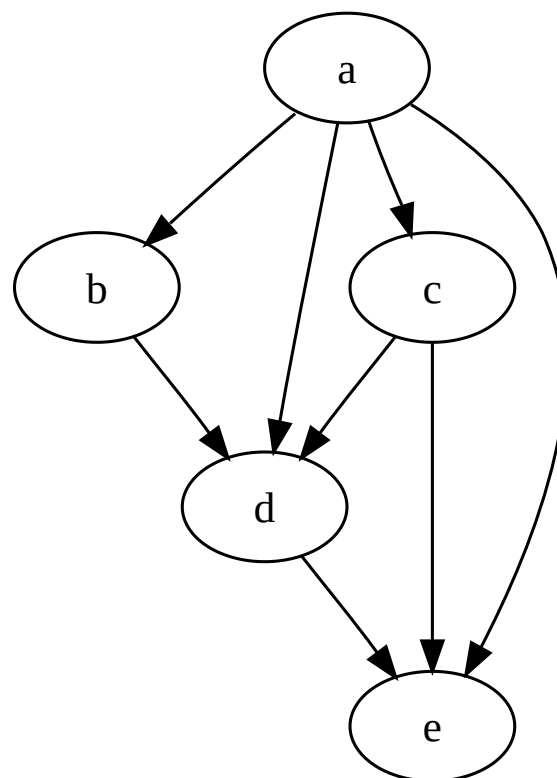


PGMs are composed of two main components: a graph and a set of probability distributions. The graph represents the structure of the relationships between the variables in the model, while the probability distributions define the probabilities of those relationships. PGMs can be divided into two main classes: directed graphical models (also known as Bayesian networks) and undirected graphical models (also known as Markov random fields).

PGMs are used in a variety of applications, such as medical diagnosis, speech recognition, and image segmentation. They allow for efficient inference and learning of complex probabilistic relationships, even in the presence of noisy or incomplete data. By modelling the uncertainty in a system, PGMs provide a powerful tool for decision-making and prediction in a wide range of domains.

## (b)

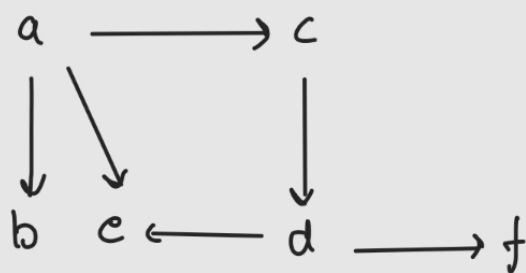
DAG stands for Directed Acyclic Graph, which is a type of graph that has directed edges between its nodes, but does not contain any directed cycles.



(c)

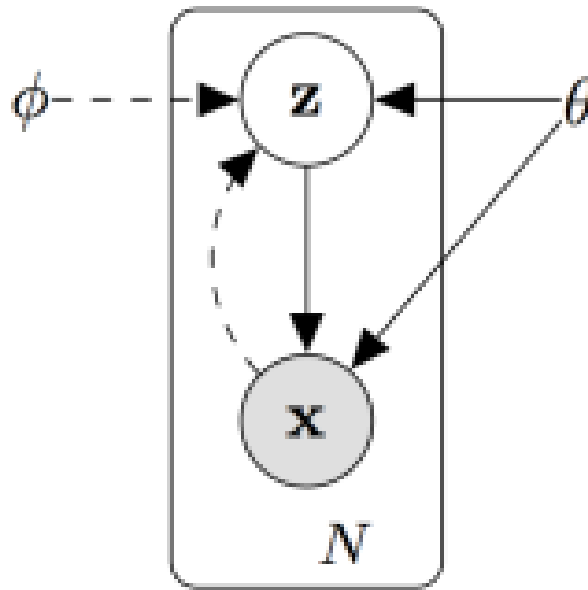
(c)

$\{a, b, c, d, e\}$



$$\begin{aligned} P(a, b, c, d, e) &= \prod_{i=1}^D P(x_i | \text{parent}(x_i)) \\ &= p(a) \cdot p(b|a) \cdot p(c|a) \cdot p(d|c) \cdot \\ &\quad p(e|d) \cdot p(f|d) \end{aligned}$$

(d)



The model used by VAE :

$$p(x, z) = p(x|z) \cdot p(z)$$

(e)

(e)

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

(f)

In VAE, importance sampling is used to estimate the expectations of the posterior distribution of the latent variables. The posterior distribution of the latent variables is intractable to compute, so importance sampling is used to obtain a lower bound that approximates it.

$$\begin{aligned}
 (f) \quad \log(p_\theta(x)) &= \log \int p_\theta(x|z) p(z) dz \\
 &\quad \searrow \text{intractable} \\
 &= \log \int q(z) \cdot \frac{p_\theta(x|z) p(z)}{q(z)} dz \quad \left[ \text{Importance Sampling} \right] \\
 &\geq \int q(z) \log \left( \frac{p_\theta(x|z) p(z)}{q(z)} \right) dz \\
 &= \underbrace{E_{q(z)} [\log p_\theta(x|z)]}_{\searrow \text{ELBO}} - D_{KL}(q(z) \parallel p(z))
 \end{aligned}$$

(g)

One common choice for the importance sampling distribution is the proposal distribution, which is typically a simpler and more tractable distribution than the true posterior. The proposal distribution is used to generate samples from the latent space, which are then used to estimate the expectations needed to compute the loss function and perform inference.

One common choice for the proposal distribution in VAE is the multivariate Gaussian distribution, which is easy to sample from and is often a good approximation to the true posterior distribution.

**(h)**

**Generalized Jensen's Inequality**

If a function  $g(x)$  maps inputs to scalar outputs in  $R$  and  $f : R \rightarrow R$  is a convex function, then for any distribution  $P_X(x)$ :

$$E_{P_X(x)}[f(g(x))] \geq f(E_{P_X(x)}[g(x)])$$

**Proof:**

(h)

Proof

Say  $y = g(x)$ ,

then by LOTUS:

$$E_{p_x(x)}[f(g(x))] = E_{p_y(y)}[f(y)]$$

$$\geq f(E_{p_y(y)}[y]) \quad \left[ \begin{array}{l} \text{By} \\ \text{Jensen's} \\ \text{inequality} \end{array} \right]$$

$$= f(E_{p_x(x)}[g(x)]) \quad \left[ \begin{array}{l} \text{Substituting} \\ g(x) \text{ back} \\ \text{by LOTUS} \end{array} \right]$$

For a concave func<sup>n</sup> 'f' the inequality is as follows:

$$E_{p_x(x)}[f(g(x))] \leq f(E_{p_x(x)}[g(x)])$$

(i)

A latent variable is a variable that is not directly observed but is inferred from other variables that are observed. Latent variables are used to represent underlying

factors or concepts that cannot be directly measured or observed, but are believed to have an effect on the observed variables.

Latent variables can be used in various models, such as factor analysis, latent class analysis, and latent variable models, to explain the relationships between observable variables and to capture the unobservable factors that influence them.

## **(j)**

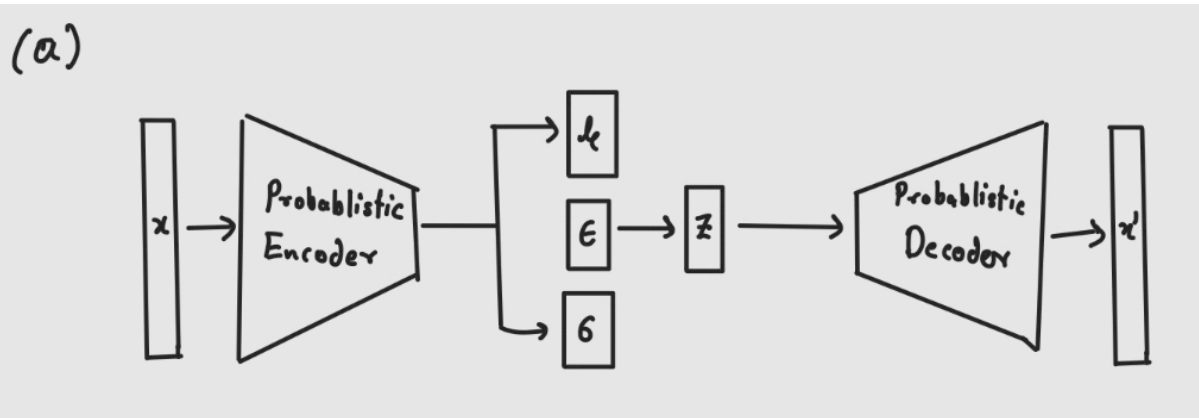
The use of a latent variable in VAE has several benefits:

1. Dimensionality reduction: By representing the data in a low-dimensional latent space, VAE can reduce the dimensionality of the data, making it easier to process and analyse.
2. Continuous representations: VAE can generate a continuous representation of the data in the latent space, allowing for smoother interpolation between different data points.
3. Generative modelling: The decoder network in VAE can be used to generate new data points by sampling from the latent space, which is represented by a distribution.
4. Regularization: By introducing a prior distribution over the latent space, VAE can regularize the latent variable, preventing overfitting and improving generalization.

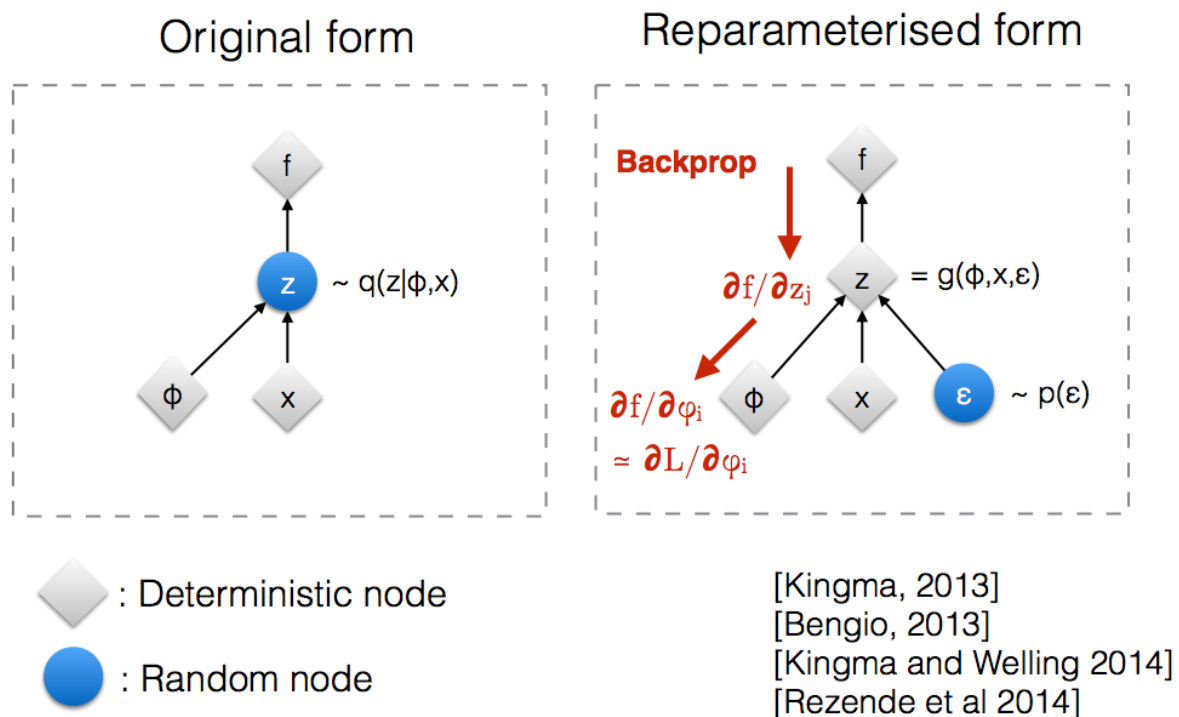
## **Problem-2**

### **(a)**





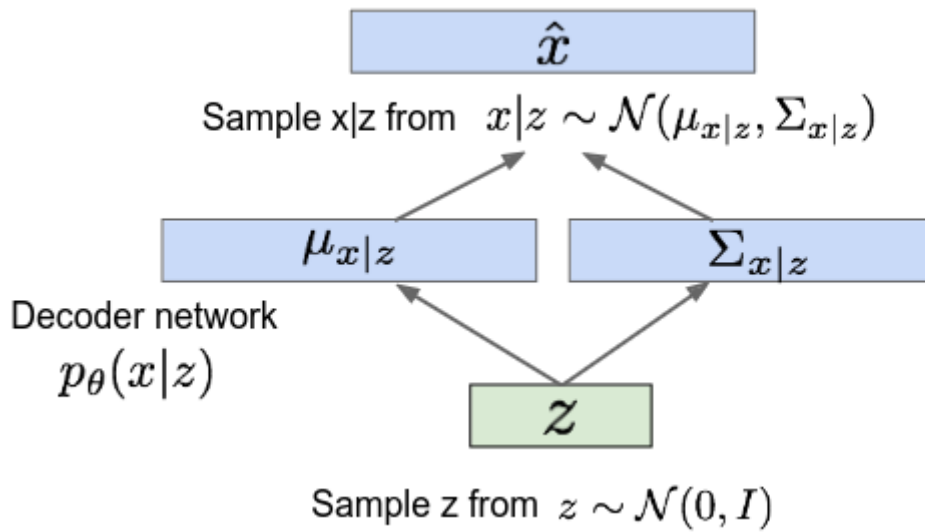
(b)



The reparametrization trick is a technique used in Variational Autoencoder (VAE) models to enable backpropagation through the stochastic layer, which is typically represented by a random sample from a normal distribution.

The reparametrization trick separates the random sampling operation from the model parameters, allowing gradients to flow through the stochastic layer and enabling efficient training via stochastic gradient descent.

(c)



Once a Variational Autoencoder (VAE) is trained, it can be used to generate new images by sampling from the learned latent space. The generation process typically involves the following steps:

1. Sampling a point from the latent space
2. Decoding the latent point back to the original input space, generating a new image.
3. Repeat the above steps to generate multiple new images by sampling different points from the learned latent space.

(d)

(d)

$$\text{Data likelihood: } p_{\theta}(n) = \int \underbrace{p_{\theta}(z)}_{\substack{\downarrow \\ \text{we want} \\ \text{to maximize} \\ \text{this}}} \underbrace{p_{\theta}(n|z)}_{\substack{\downarrow \\ \text{intractable}}} dz$$

$\Rightarrow$  Using variational inference to approximate the unknown posterior distribution from only the observed data.

$$\log(p_{\theta}(n)) = E_{z \sim q_{\phi}(z|n)} [\log p_{\theta}(n)] \quad (p_{\theta}(n) \text{ doesn't depend on } z)$$

$$= E_z \left[ \log \left( \frac{p_{\theta}(n|z) p_{\theta}(z)}{p_{\theta}(z|n)} \right) \right] \quad (\text{Bayes})$$

$$= E_z \left[ \log \left( \frac{p_{\theta}(n|z) p_{\theta}(z)}{p_{\theta}(z|n)} \cdot \frac{q_{\phi}(z|n)}{q_{\phi}(z|n)} \right) \right] \quad (\text{importance sampling})$$

$$= E_z [\log(p_{\theta}(n|z))] - E_z \left[ \log \left( \frac{q_{\phi}(z|n)}{p_{\theta}(z)} \right) \right]$$

$$+ E_z \left[ \log \left( \frac{q_{\phi}(z|n)}{p_{\theta}(z|n)} \right) \right]$$

$$= E_z [\log p_{\theta}(n|z)] - D_{KL}[q_{\phi}(z|n) \| p_{\theta}(z)] + \underbrace{D_{KL}[q_{\phi}(z|n) \| p_{\theta}(z|n)]}_{\substack{\uparrow \\ \text{intractable}}} \geq 0$$

$$\geq \underbrace{E_z [\log p_{\theta}(n|z)] - D_{KL}[q_{\phi}(z|n) \| p_{\theta}(z)]}_{\mathcal{L}(n, \theta, \phi)}$$

$\downarrow$   
Tractable Lower Bound  
(ELBO)

### **(e)**

Monte Carlo estimation is used in two places in a Variational Autoencoder (VAE):

1. Computation of ELBO: To compute the ELBO, Monte Carlo integration is used to estimate the expected log-likelihood under the approximate posterior distribution over the latent variables. Specifically, a set of samples is drawn from the posterior distribution, and the log-likelihood of the data is estimated by averaging the log-likelihood over the samples.
2. Generating of new samples: Once the VAE is trained, it can be used to generate new samples by sampling from the learned latent space. To do this, Monte Carlo sampling is used to generate a set of samples from the prior distribution over the latent variables.

### **(f)**

(f)

Generalize to 2 multivariate Normal Distributions

$$p(x) = N(x; \mu_1, \Sigma_1)$$

$\mu_1, \mu_2 \rightarrow$  Means

— given

$$q(x) = N(x; \mu_2, \Sigma_2)$$

$\Sigma_1, \Sigma_2 \rightarrow$  Covariance Matrices

$$D_{KL}(p(x) \| q(x)) = \left\{ \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right] \right\}$$

Proof

$$KL(p \| q) = \sum_n p \log \left( \frac{p}{q} \right)$$

$$= \sum_n p (\log p - \log q)$$

$$\left\{ \begin{array}{l} \text{Multivariate Normal Density} \\ N(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right) \\ \downarrow \\ \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \end{array} \right\}$$

(

$$KL(p \| q) = \sum_n p(x) \left\{ \frac{1}{2} \log \left( \frac{|\Sigma_2|}{|\Sigma_1|} \right) - \frac{1}{2} \underbrace{(\mu_2 - \mu_1)^T \Sigma_1^{-1} (\mu_2 - \mu_1)} + \frac{1}{2} (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right\}$$

## Part - 1

$$= \sum_n p(n) \left[ \frac{1}{2} (n - p_1)^T \Sigma_1^{-1} (n - p_1) \right]$$

$$\Rightarrow = E \left[ \frac{1}{2} (n - p_1)^T \Sigma_1^{-1} (n - p_1) \right]$$

properties of trace

$E(x) = E[\text{tr}(x)]$ , if  $x$ -scalar    using  $\Rightarrow$

$\text{tr}(AB) = \text{tr}(BA)$

$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$

$E[\text{tr}(A)] = \text{tr}(E[A])$

$$\begin{aligned} E(x^T A x) &= E[\text{tr}(x^T A x)] \\ &= E[\text{tr}(A x x^T)] \\ &= \text{tr}(E[A x x^T]) \end{aligned}$$

$$\text{tr} \left( E_p \left[ \frac{1}{2} (n - p_1)^T \Sigma_1^{-1} (n - p_1) \right] \right)$$

$$= \text{tr} \left( E_p \left[ \frac{1}{2} \underbrace{(n - p_1)(n - p_1)^T}_{\text{covariance}} \Sigma_1^{-1} \right] \right)$$

$$= \text{tr} \left( E_p \left[ \underbrace{(n - p_1)(n - p_1)^T}_{\text{covariance}} \right] \frac{1}{2} \Sigma_1^{-1} \right)$$

$$= \text{tr} \left( \frac{1}{2} \Sigma_1 \Sigma_1^{-1} \right)$$

$$= \text{tr}(I_K)$$

$$= \boxed{K} \longrightarrow \propto \text{constant}$$

## Part-2

$$\sum_n p(n) \left[ \frac{1}{2} (n - p_2)^T \Sigma_2^{-1} (n - p_2) \right]$$

$$= \sum_n p(n) \left[ \frac{1}{2} ((n - p_1) + (p_1 - p_2))^T \Sigma_2^{-1} ((n - p_1) + (p_1 - p_2)) \right]$$

↓

$$\left. \begin{aligned} & (A+B)^T \Sigma_2^{-1} (A+B) \\ \Rightarrow & (A^T + B^T) \Sigma_2^{-1} (A+B) \\ \Rightarrow & A^T \Sigma_2^{-1} A + B^T \Sigma_2^{-1} B + \underbrace{(B^T \Sigma_2^{-1} A + A^T \Sigma_2^{-1} B)}_{\text{same}} \end{aligned} \right\}$$

⋮

↖ similar to prev.

$$\Rightarrow E \left[ \frac{1}{2} (n - p_1)^T \Sigma_2^{-1} (n - p_1) \right] + E \left[ (n - p_1)^T \Sigma_2^{-1} (p_1 - p_2) \right]$$

$$+ E \left[ (p_1 - p_2)^T \Sigma_2^{-1} (p_1 - p_2) \right]$$

↓  
constant

0  
(proof?)

$$\Rightarrow \left\{ \frac{\Sigma_2^{-1} \Sigma_1}{2} \right\} + \underbrace{(p_1 - p_2)^T \Sigma_2^{-1} (p_1 - p_2)}_{\substack{\downarrow \\ E[\text{const.}] = \text{const.} \\ \downarrow \\ \beta_{\text{const.}}}} + 0$$

↓  
 $E[\text{const.}] = \text{const.}$

↓  
 $\beta_{\text{const.}}$

$$\left\{ \begin{aligned} & E_p [(n - p_1)^T \Sigma_2^{-1} (p_2 - p_1)] \\ &= [(E_p[n] - p_1)^T \Sigma_2^{-1} (p_2 - p_1)] \\ &= [(p_1 - p_1)^T \dots] = \boxed{0} \end{aligned} \right\}$$

## (g)

Autoencoder:

- Non-probabilistic model
- No assumption about the data distribution
- Encoder maps input to a fixed-length latent representation
- No explicit way to sample from the learned latent space

Variational Autoencoder (VAE):

- Probabilistic model
- Assumes a prior distribution over the latent space
- Encoder maps input to a distribution over the latent representation
- Allows for explicit sampling from the learned latent space