

Variational Auto-Encoders

Wilker Aziz

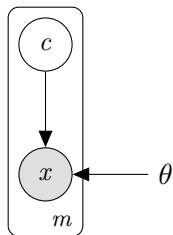
Universiteit van Amsterdam

`w.aziz@uva.nl`

May 19, 2017

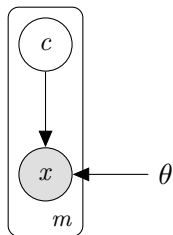
Generative models with neural networks

Mixture model



Generative models with neural networks

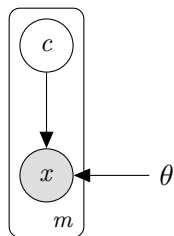
Mixture model



- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$

Generative models with neural networks

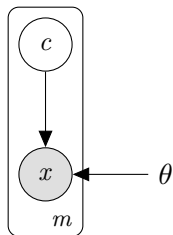
Mixture model



- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$

Generative models with neural networks

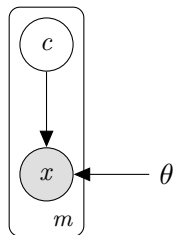
Mixture model



- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$
- ▶ where $P(X|C = c) = \text{Cat}(f_{\theta}(c))$

Generative models with neural networks

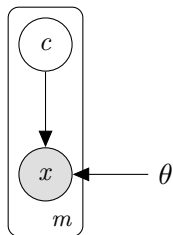
Mixture model



- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$
- ▶ where $P(X|C = c) = \text{Cat}(f_{\theta}(c))$
 - ▶ e.g. $f_{\theta}(c) = \text{softmax}(W^{(f)}g(c) + b^{(f)})$

Generative models with neural networks

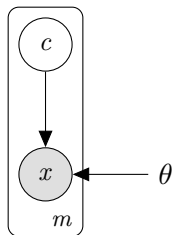
Mixture model



- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$
- ▶ where $P(X|C = c) = \text{Cat}(f_{\theta}(c))$
 - ▶ e.g. $f_{\theta}(c) = \text{softmax}(W^{(f)}g(c) + b^{(f)})$
and $g(c) = \tanh(W^{(g)}r(c) + b^{(g)})$

Generative models with neural networks

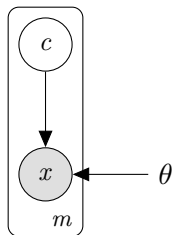
Mixture model



- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$
- ▶ where $P(X|C = c) = \text{Cat}(f_{\theta}(c))$
 - ▶ e.g. $f_{\theta}(c) = \text{softmax}(W^{(f)}g(c) + b^{(f)})$
and $g(c) = \tanh(W^{(g)}r(c) + b^{(g)})$
and $r(c) = Ec$

Generative models with neural networks

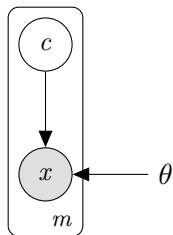
Mixture model



- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$
- ▶ where $P(X|C = c) = \text{Cat}(f_{\theta}(c))$
 - ▶ e.g. $f_{\theta}(c) = \text{softmax}(W^{(f)}g(c) + b^{(f)})$
and $g(c) = \tanh(W^{(g)}r(c) + b^{(g)})$
and $r(c) = Ec$
- ▶ with $\theta = (E, W^{(f)}, b^{(f)}, W^{(g)}, b^{(g)})$

Generative models with neural networks

Mixture model

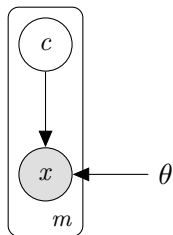


- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$
- ▶ where $P(X|C = c) = \text{Cat}(f_{\theta}(c))$
 - ▶ e.g. $f_{\theta}(c) = \text{softmax}(W^{(f)}g(c) + b^{(f)})$
and $g(c) = \tanh(W^{(g)}r(c) + b^{(g)})$
and $r(c) = Ec$
- ▶ with $\theta = (E, W^{(f)}, b^{(f)}, W^{(g)}, b^{(g)})$

$$P(x) = \sum_{c=1}^K \underbrace{P(c)P(x|c)}_{\substack{\text{differentiable function of } \theta \\ \text{tractable for small } K}}$$

Generative models with neural networks

Mixture model



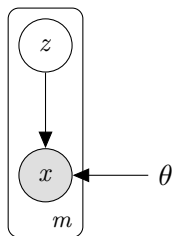
- ▶ sample a latent class $c \in \{1, \dots, K\}$
 $c \sim U(\frac{1}{K})$
- ▶ generate categorical observation x from c
 $x \sim P(X|C = c)$
- ▶ where $P(X|C = c) = \text{Cat}(f_\theta(c))$
 - ▶ e.g. $f_\theta(c) = \text{softmax}(W^{(f)}g(c) + b^{(f)})$
and $g(c) = \tanh(W^{(g)}r(c) + b^{(g)})$
and $r(c) = Ec$
- ▶ with $\theta = (E, W^{(f)}, b^{(f)}, W^{(g)}, b^{(g)})$

$$P(x) = \sum_{c=1}^K \underbrace{P(c)P(x|c)}_{\substack{\text{differentiable function of } \theta \\ \text{tractable for small } K}}$$

Gradient-based optimisation! $\nabla_\theta \log P_\theta(x)$

Generative models with neural networks

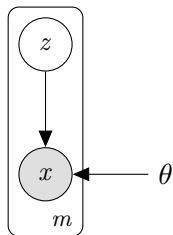
Continuous mixture model



Generative models with neural networks

Continuous mixture model

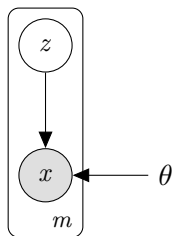
- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$



Generative models with neural networks

Continuous mixture model

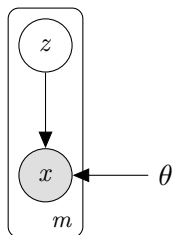
- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$
- ▶ generate categorical observation x from z
 $x \sim P(X|Z = z)$



Generative models with neural networks

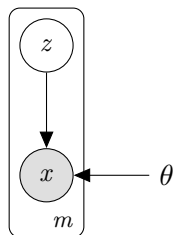
Continuous mixture model

- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$
- ▶ generate categorical observation x from z
 $x \sim P(X|Z = z)$
- ▶ where $P(X|Z = z) = \text{Cat}(f_\theta(z))$



Generative models with neural networks

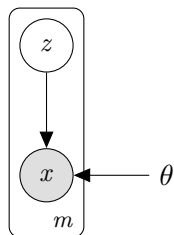
Continuous mixture model



- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$
- ▶ generate categorical observation x from z
 $x \sim P(X|Z = z)$
- ▶ where $P(X|Z = z) = \text{Cat}(f_\theta(z))$
 - ▶ e.g. $f_\theta(z) = \text{softmax}(W^{(f)}g(z) + b^{(f)})$

Generative models with neural networks

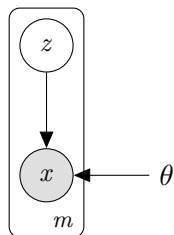
Continuous mixture model



- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$
- ▶ generate categorical observation x from z
 $x \sim P(X|Z = z)$
- ▶ where $P(X|Z = z) = \text{Cat}(f_\theta(z))$
 - ▶ e.g. $f_\theta(z) = \text{softmax}(W^{(f)}g(z) + b^{(f)})$
and $g(z) = \tanh(W^{(g)}z + b^{(g)})$

Generative models with neural networks

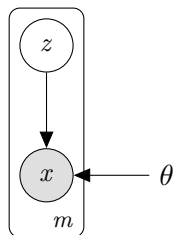
Continuous mixture model



- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$
- ▶ generate categorical observation x from z
 $x \sim P(X|Z = z)$
- ▶ where $P(X|Z = z) = \text{Cat}(f_\theta(z))$
 - ▶ e.g. $f_\theta(z) = \text{softmax}(W^{(f)}g(z) + b^{(f)})$
and $g(z) = \tanh(W^{(g)}z + b^{(g)})$
- ▶ with $\theta = (W^{(f)}, b^{(f)}, W^{(g)}, b^{(g)})$

Generative models with neural networks

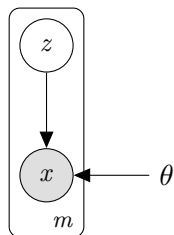
Continuous mixture model



- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$
- ▶ generate categorical observation x from z
 $x \sim P(X|Z = z)$
- ▶ where $P(X|Z = z) = \text{Cat}(f_\theta(z))$
 - ▶ e.g. $f_\theta(z) = \text{softmax}(W^{(f)}g(z) + b^{(f)})$
and $g(z) = \tanh(W^{(g)}z + b^{(g)})$
- ▶ with $\theta = (W^{(f)}, b^{(f)}, W^{(g)}, b^{(g)})$
- ▶ **Intractability**

Generative models with neural networks

Continuous mixture model



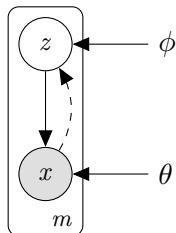
- ▶ sample a latent embedding $z \in \mathbb{R}^d$
 $z \sim \mathcal{N}(0, I)$
- ▶ generate categorical observation x from z
 $x \sim P(X|Z = z)$
- ▶ where $P(X|Z = z) = \text{Cat}(f_\theta(z))$
 - ▶ e.g. $f_\theta(z) = \text{softmax}(W^{(f)}g(z) + b^{(f)})$
and $g(z) = \tanh(W^{(g)}z + b^{(g)})$
- ▶ with $\theta = (W^{(f)}, b^{(f)}, W^{(g)}, b^{(g)})$
- ▶ **Intractability**
 - ▶ $P(x) = \int p(z)P(x|z)dz$
 - ▶ $P(z|x) = \frac{p(z)P(x|z)}{\int p(z')P(x|z')dz'}$

Variational inference

but we know VI :D

Variational inference

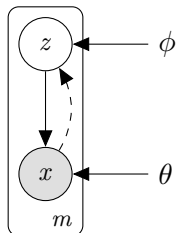
but we know VI :D



- approximate the posterior with $q_{\phi}(Z|x) = \mathcal{N}(\mu_{\phi}(x), I\sigma_{\phi}^2(x))$

Variational inference

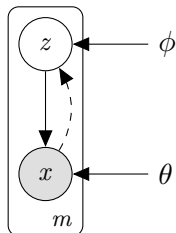
but we know VI :D



- ▶ approximate the posterior with
$$q_{\phi}(Z|x) = \mathcal{N}(\mu_{\phi}(x), I\sigma_{\phi}^2(x))$$
- ▶ where
 - ▶ $\mu_{\phi}(x) = W^{(\mu)}u(x) + b^{(\mu)}$
e.g. $u(x) = \tanh(W^{(u)}r(x) + b^{(u)})$
and $r(x) = E^{(u)}x$

Variational inference

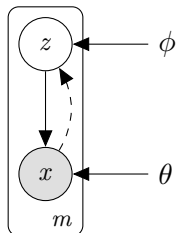
but we know VI :D



- ▶ approximate the posterior with
$$q_{\phi}(Z|x) = \mathcal{N}(\mu_{\phi}(x), I\sigma_{\phi}^2(x))$$
- ▶ where
 - ▶ $\mu_{\phi}(x) = W^{(\mu)}u(x) + b^{(\mu)}$
e.g. $u(x) = \tanh(W^{(u)}r(x) + b^{(u)})$
and $r(x) = E^{(u)}x$
 - ▶ $\sigma_{\phi}^2(x) = \exp(W^{(\sigma)}v(x) + b^{(\sigma)})$
e.g. $v(x) = \tanh(W^{(v)}r(x) + b^{(v)})$
and $r(x) = E^{(v)}x$

Variational inference

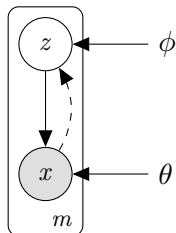
but we know VI :D



- ▶ approximate the posterior with
$$q_{\phi}(Z|x) = \mathcal{N}(\mu_{\phi}(x), I\sigma_{\phi}^2(x))$$
- ▶ where
 - ▶ $\mu_{\phi}(x) = W^{(\mu)}u(x) + b^{(\mu)}$
e.g. $u(x) = \tanh(W^{(u)}r(x) + b^{(u)})$
and $r(x) = E^{(u)}x$
 - ▶ $\sigma_{\phi}^2(x) = \exp(W^{(\sigma)}v(x) + b^{(\sigma)})$
e.g. $v(x) = \tanh(W^{(v)}r(x) + b^{(v)})$
and $r(x) = E^{(v)}x$
- ▶ with $\phi = (E^{(u,v)}, W^{(u,v,\mu,\sigma)}, b^{(u,v,\mu,\sigma)})$

Variational inference

but we know VI :D



- ▶ approximate the posterior with
$$q_{\phi}(Z|x) = \mathcal{N}(\mu_{\phi}(x), I\sigma_{\phi}^2(x))$$
- ▶ where
 - ▶ $\mu_{\phi}(x) = W^{(\mu)}u(x) + b^{(\mu)}$
e.g. $u(x) = \tanh(W^{(u)}r(x) + b^{(u)})$
and $r(x) = E^{(u)}x$
 - ▶ $\sigma_{\phi}^2(x) = \exp(W^{(\sigma)}v(x) + b^{(\sigma)})$
e.g. $v(x) = \tanh(W^{(v)}r(x) + b^{(v)})$
and $r(x) = E^{(v)}x$
- ▶ with $\phi = (E^{(u,v)}, W^{(u,v,\mu,\sigma)}, b^{(u,v,\mu,\sigma)})$

Mean field assumption

- ▶ $q_{\phi_i}(Z|x_i)$ is specified for each observation x_i by locally predicting its mean and variance

Approximate inference by optimisation

Maximise ELBO

$$\log P_{\theta}(x) \geq \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} \left[\log \frac{q_{\phi}(Z|x)}{p_{\theta}(Z)} \right]}_{-\text{KL}(q_{\theta}(Z|x)||p_{\theta}(Z))} + \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} [\log P_{\theta}(X = x|Z)]}_{\text{intractable!}}$$

Approximate inference by optimisation

Maximise ELBO

$$\log P_{\theta}(x) \geq \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} \left[\log \frac{q_{\phi}(Z|x)}{p_{\theta}(Z)} \right]}_{-\text{KL}(q_{\theta}(Z|x)||p_{\theta}(Z))} + \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} [\log P_{\theta}(X = x|Z)]}_{\text{intractable!}}$$

Prior term

$$\text{KL}(q_{\phi}(Z|x)||p_{\theta}(Z)) = -\frac{1}{2} \sum_{j=1}^d (1 + \log \sigma_{\phi}^2(x)_j - \mu_{\phi}^2(x)_j - \sigma_{\phi}^2(x)_j)$$

Approximate inference by optimisation

Maximise ELBO

$$\log P_{\theta}(x) \geq \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} \left[\log \frac{q_{\phi}(Z|x)}{p_{\theta}(Z)} \right]}_{-\text{KL}(q_{\phi}(Z|x)||p_{\theta}(Z))} + \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} [\log P_{\theta}(X = x|Z)]}_{\text{intractable!}}$$

Prior term

$$\text{KL}(q_{\phi}(Z|x)||p_{\theta}(Z)) = -\frac{1}{2} \sum_{j=1}^d (1 + \log \sigma_{\phi}^2(x)_j - \mu_{\phi}^2(x)_j - \sigma_{\phi}^2(x)_j)$$

Likelihood term is intractable

- ▶ the Categorical likelihood is not conjugate with the Normal approximate posterior

Change of variable for location-scale distributions

For $Z \sim \mathcal{N}(\mu, \sigma^2)$ we can re-express Z in terms of $E \sim \mathcal{N}(0, 1)$

- ▶ $Z = \mu + \sigma E$

Change of variable for location-scale distributions

For $Z \sim \mathcal{N}(\mu, \sigma^2)$ we can re-express Z in terms of $E \sim \mathcal{N}(0, 1)$

► $Z = \mu + \sigma E$

then we can re-express expectations

$$\mathbb{E}_{\mathcal{N}(\mu, \sigma^2)}[f(Z)] = \mathbb{E}_{\mathcal{N}(0, I)}[f(\mu + \sigma E)]$$

Change of variable for location-scale distributions

For $Z \sim \mathcal{N}(\mu, \sigma^2)$ we can re-express Z in terms of $E \sim \mathcal{N}(0, 1)$

► $Z = \mu + \sigma E$

then we can re-express expectations

$$\mathbb{E}_{\mathcal{N}(\mu, \sigma^2)}[f(Z)] = \mathbb{E}_{\mathcal{N}(0, I)}[f(\mu + \sigma E)]$$

back to the ELBO

$$\mathbb{E}_{q_\phi(Z|x)} [\log P(x|Z)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\log P(x|Z = \mu_\phi(x) + \sigma_\phi(x)\epsilon)]$$

Monte Carlo estimate

$$\begin{aligned}\mathbb{E}_{q_\phi(Z|x)} [\log P(x|Z)] &= \mathbb{E}_{\epsilon \sim N(0, I)} [\log P(x|Z = \mu_\phi(x) + \sigma_\phi(x)\epsilon)] \\ &\approx \frac{1}{N} \sum_{n=1}^N \log P\left(x|\mu_\phi(x) + \sigma_\phi(x)\epsilon^{(n)}\right)\end{aligned}$$

MC estimate of the ELBO

$$\begin{aligned}\log P_{\theta}(x) &\geq \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} \left[\log \frac{q_{\phi}(Z|x)}{p_{\theta}(Z)} \right]}_{-\text{KL}(q_{\theta}(Z|x)||p_{\theta}(Z))} + \underbrace{\mathbb{E}_{q_{\phi}(Z|x)} [\log P_{\theta}(X = x|Z)]}_{\text{intractable!}} \\ &\approx \underbrace{\frac{1}{2} \sum_{j=1}^d \left(1 + \log \sigma_{\phi}^2(x)_j - \mu_{\phi}^2(x)_j - \sigma_{\phi}^2(x)_j \right)}_{-\text{KL}(q_{\theta}(Z|x)||p_{\theta}(Z))} \\ &\quad + \underbrace{\log P_{\theta}(x|\mu_{\phi}(x) + \sigma_{\phi}(x)\epsilon)}_{\text{single-sample estimate}}\end{aligned}$$

Gradient-based optimisation

Let $\mathcal{L}(\theta, \phi|x)$ be our objective function

$$\mathcal{L}(\theta, \phi|x) = \underbrace{\frac{1}{2} \sum_{j=1}^d \left(1 + \log \sigma_{\phi}^2(x)_j - \mu_{\phi}^2(x)_j - \sigma_{\phi}^2(x)_j \right)}_{\text{differentiable function of } \phi}$$
$$+ \underbrace{\log P_{\theta}(x|\mu_{\phi}(x) + \sigma_{\phi}(x)\epsilon)}_{\text{differentiable function of } \theta \text{ and } \phi}$$

Gradient-based optimisation

Let $\mathcal{L}(\theta, \phi|x)$ be our objective function

$$\mathcal{L}(\theta, \phi|x) = \underbrace{\frac{1}{2} \sum_{j=1}^d \left(1 + \log \sigma_{\phi}^2(x)_j - \mu_{\phi}^2(x)_j - \sigma_{\phi}^2(x)_j \right)}_{\text{differentiable function of } \phi} + \underbrace{\log P_{\theta}(x|\mu_{\phi}(x) + \sigma_{\phi}(x)\epsilon)}_{\text{differentiable function of } \theta \text{ and } \phi}$$

We can update θ and ϕ using stochastic gradient steps

- ▶ we know chain rule (thus we can get a gradient)
- ▶ we have a noisy though unbiased estimate
- ▶ guaranteed convergence to a local optimum of \mathcal{L} (with appropriate learning rate schedule)

Further reading

- ▶ Auto-Encoding variational Bayes [Kingma and Welling, 2014]

References I

Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.