

Approximate Inference in Latent Variable Models based on NNs

Pablo M. Olmos, olmos@tsc.uc3m.es

October 3, 2017

uc3m | Universidad **Carlos III** de Madrid



- Latent Variable Models
- Variational Inference
- Neural Networks
- Stochastic Gradient Descent

NIPS 2016 word cloud

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- VAE Objective Function
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

- Normalizing flows
- Adversarial networks and Implicit Variational Inference

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- VAE Objective Function
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

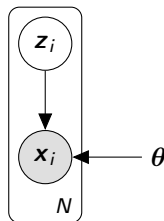
- Normalizing flows
- Adversarial networks and Implicit Variational Inference

Fitting a distribution using VAEs: Generative Model

Consider a set of i.i.d observations $\mathbf{x}^{(i)} \in \mathbb{R}^d$, $i = 1, \dots, N$. Let $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$.

Generative Model using a Latent Space

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

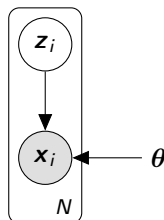


Fitting a distribution using VAEs: Generative Model

Consider a set of i.i.d observations $\mathbf{x}^{(i)} \in \mathbb{R}^d$, $i = 1, \dots, N$. Let $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$.

Generative Model using a Latent Space

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$



- $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$ ($k \leq d$) (low-dimensional embedding)
- $\mathbf{x}|\mathbf{z} \sim \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$
- $\mu_{\theta}(\mathbf{z})$ and $\log(\sigma_{\theta}(\mathbf{z}))$ are the outputs of a NN $\mathbb{R}^k \rightarrow \mathbb{R}^d$ with parameter vector θ (Decoding network).

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- **VAE Objective Function**
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

- Normalizing flows
- Adversarial networks and Implicit Variational Inference

Fitting a distribution using VAEs: ML objective

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

ML estimation (in general intractable)

$$\begin{aligned}\theta^* &= \arg \max_{\theta \in \Omega} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) \\ &\approx \arg \max_{\theta \in \Omega} \int p^*(\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{x} \\ &= \arg \min_{\theta \in \Omega} \text{KL}(p^*(\mathbf{x}) || p_{\theta}(\mathbf{x}))\end{aligned}$$

Recall that $\text{KL}(p(\mathbf{x}) || q(\mathbf{x})) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \geq 0$

Fitting a distribution using VAEs: Variational objective

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

Consider **any** approximation $q(\mathbf{z})$ to $p_{\theta}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$

Fitting a distribution using VAEs: Variational objective

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

Consider **any** approximation $q(\mathbf{z})$ to $p_{\theta}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$

$$\log p_{\theta}(\mathbf{x}) = \text{KL}(q(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x})) + \mathcal{L}(\mathbf{x}, \theta)$$

where

$$\mathcal{L}(\mathbf{x}, \theta) = \int q(\mathbf{z}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z} = \mathbb{E}_{q(\mathbf{z})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q(\mathbf{z})||p(\mathbf{z}))$$

is the Evidence Lower Bound (ELBO).

We will optimize $\mathcal{L}(\mathbf{x}, \theta)$ w.r.t. θ But first we need to select a variational family for $q(\mathbf{z})$.

Fitting a distribution using VAEs: Variational objective (II)

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

The Inference Network or Encoding Network

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta,\mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

where $\mu_{\eta}(\mathbf{x})$ and $\log(\sigma_{\eta}(\mathbf{x}))$ are the outputs of a NN $\mathbb{R}^d \rightarrow \mathbb{R}^k$ with parameter vector η .

Fitting a distribution using VAEs: Variational objective (II)

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

The Inference Network or Encoding Network

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta,\mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

where $\mu_{\eta}(\mathbf{x})$ and $\log(\sigma_{\eta}(\mathbf{x}))$ are the outputs of a NN $\mathbb{R}^d \rightarrow \mathbb{R}^k$ with parameter vector η .

The VAE optimization problem

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \text{KL} \left(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z}) \right) \right)$$

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- VAE Objective Function
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

- Normalizing flows
- Adversarial networks and Implicit Variational Inference

Fitting a distribution using VAEs: Stochastic Optimization

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

Recall that $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$, thus

$$\text{KL} (q_{\eta, \mathbf{x}}(\mathbf{z}) || p(\mathbf{z})) = \frac{1}{2} \left(-k + \sum_{j=1}^k \sigma_{\eta, j}(\mathbf{x}) - \log \sigma_{\eta, j}(\mathbf{x}) + \mu_{\eta, j}^2 \right)$$

→ Differentiable w.r.t. η

Fitting a distribution using VAEs: Stochastic Optimization

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

Recall that $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$, thus

$$\text{KL} (q_{\eta, \mathbf{x}}(\mathbf{z}) || p(\mathbf{z})) = \frac{1}{2} \left(-k + \sum_{j=1}^k \sigma_{\eta, j}(\mathbf{x}) - \log \sigma_{\eta, j}(\mathbf{x}) + \mu_{\eta, j}^2 \right)$$

→ Differentiable w.r.t. η

Unbiased gradient estimator

$$\nabla_{\eta} \left(\frac{1}{N} \sum_{i=1}^N \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right) \approx \nabla_{\eta} \left(\frac{1}{M} \sum_{i \in \mathcal{M}} \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

where \mathcal{M} represents a M -sized minibatch of data sampled at random from \mathcal{D} ($M \ll N$).

Fitting a distribution using VAEs: Stochastic Optimization (II)

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

We also need unbiased gradient estimates for

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right]$$

Fitting a distribution using VAEs: Stochastic Optimization (II)

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

We also need unbiased gradient estimates for

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right]$$

We use a Monte Carlo sampling estimator

$$\mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] \approx \frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{i,s})$$

where $\mathbf{z}^{(s,i)}$, $s = 1, \dots, S$ are i.i.d. samples from $q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})$.

We typically use a **single** sample, i.e., $S = 1$ (huge estimator variance, but cheap computation).

Fitting a distribution using VAEs: Stochastic Optimization (III)

If $\mathbf{z}^{(s,i)}$ are i.i.d. samples from $q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}^{(i)}), \text{diag}(\sigma_{\eta}(\mathbf{x}^{(i)})))$

$$\frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{i,s}) \rightarrow \text{how do we compute gradients w.r.t. } \boldsymbol{\eta}?$$

Fitting a distribution using VAEs: Stochastic Optimization (III)

If $\mathbf{z}^{(s,i)}$ are i.i.d. samples from $q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}^{(i)}), \text{diag}(\sigma_{\eta}(\mathbf{x}^{(i)})))$

$$\frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{i,s}) \rightarrow \text{how do we compute gradients w.r.t. } \eta?$$

Reparameterization Trick

Express each sample $\mathbf{z}^{(s,i)}$ as a deterministic function of $\mathbf{x}^{(i)}$ and some noise vector $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ that is independent of η . For a Gaussian distribution we have

$$\mathbf{z}^{(s,i)} = f_{\eta}(\mathbf{x}^{(i)}, \epsilon) = \mu_{\eta}(\mathbf{x}^{(i)}) + \sqrt{\sigma_{\eta}(\mathbf{x}^{(i)})} \cdot \epsilon$$

Fitting a distribution using VAEs: Stochastic Optimization (III)

If $\mathbf{z}^{(s,i)}$ are i.i.d. samples from $q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}^{(i)}), \text{diag}(\sigma_{\eta}(\mathbf{x}^{(i)})))$

$$\frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{s,i}) \rightarrow \text{how do we compute gradients w.r.t. } \eta?$$

Reparameterization Trick

Express each sample $\mathbf{z}^{(s,i)}$ as a deterministic function of $\mathbf{x}^{(i)}$ and some noise vector $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ that is independent of η . For a Gaussian distribution we have

$$\mathbf{z}^{(s,i)} = f_{\eta}(\mathbf{x}^{(i)}, \epsilon) = \mu_{\eta}(\mathbf{x}^{(i)}) + \sqrt{\sigma_{\eta}(\mathbf{x}^{(i)})} \cdot \epsilon$$

Putting all together...

$$\mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z})] = \mathbb{E}_{\epsilon} [\log p_{\theta}(\mathbf{x}^{(i)} | f_{\eta}(\epsilon, \mathbf{x}^{(i)}))] \approx \frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x}^{(i)} | f_{\eta}(\epsilon^{s,i}, \mathbf{x}^{(i)}))$$

where $\epsilon^{s,i}$, $s = 1, \dots, S$ are i.i.d. samples from $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

Fitting a distribution using VAEs: the algorithm

Algorithm 1 The Variational Autoencoder (S=1)

- 1: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}_0, \boldsymbol{\eta} \leftarrow \boldsymbol{\eta}_0$
- 2: $\ell \leftarrow 0$
- 3: **while** not converged **do**
- 4: Sample minibatch $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$ from \mathcal{D}
- 5: Sample $\{\epsilon^1, \dots, \epsilon^M\}$ from $p(\epsilon) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6: Compute noisy gradients:

$$\mathbf{g}_{\boldsymbol{\theta}}, \mathbf{g}_{\boldsymbol{\eta}} \leftarrow \frac{1}{M} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}, \boldsymbol{\eta}} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | f_{\boldsymbol{\eta}}(\epsilon^i, \mathbf{x}^{(i)})) - \text{KL}(q_{\boldsymbol{\eta}, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right]$$

- 7: Perform SGD-updates:

$$\boldsymbol{\theta}_{\ell+1} \leftarrow \boldsymbol{\theta}_{\ell+1} + h_{\ell} \mathbf{g}_{\boldsymbol{\theta}}$$

$$\boldsymbol{\eta}_{\ell+1} \leftarrow \boldsymbol{\eta}_{\ell+1} + h_{\ell} \mathbf{g}_{\boldsymbol{\eta}}$$

- 8: $\ell \leftarrow \ell + 1$
 - 9: **end while**
-

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- VAE Objective Function
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

- Normalizing flows
- Adversarial networks and Implicit Variational Inference

- Approximate Inference with Amortised MCMC
<https://arxiv.org/pdf/1702.08343.pdf>
- Autoencoding Variational Inference for Topic Models
<https://arxiv.org/pdf/1703.01488.pdf>
- Improved Variational Autoencoders for Text Modeling using Dilated Convolutions
<http://proceedings.mlr.press/v70/yang17d/yang17d.pdf>
- Variational Sequential Monte Carlo
<https://arxiv.org/pdf/1705.11140.pdf>
- Stick Breaking Variational Autoencoders
<https://arxiv.org/pdf/1605.06197.pdf>
- AutoGP: Exploring the Capabilities and Limitations of Gaussian Process Models
<http://arxiv.org/abs/1610.05392>

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- VAE Objective Function
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

- Normalizing flows
- Adversarial networks and Implicit Variational Inference

Revisiting the VAE generative model

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \text{KL} \left(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z}) \right) \right)$$

Revisiting the VAE generative model

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta,\mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \text{KL} \left(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z}) \right) \right)$$

Some critical assumptions (among many others)

- Unimodal posterior approximation.
- $p_{\theta}(\mathbf{x}|\mathbf{z})$ is known.
- I can find a valid reconstruction $p_{\theta}(\mathbf{x}|\mathbf{z})$ model.
- Prior too simple? It does not enforce interpretability in the latent space.
- ...

Revisiting the VAE generative model

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta,\mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \text{KL} \left(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z}) \right) \right)$$

Some critical assumptions (among many others)

- **Unimodal posterior approximation.**
- $p_{\theta}(\mathbf{x}|\mathbf{z})$ is known.
- I can find a valid reconstruction $p_{\theta}(\mathbf{x}|\mathbf{z})$ model.
- Prior too simple? It does not enforce interpretability in the latent space.
- ...

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- VAE Objective Function
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

- Normalizing flows
- Adversarial networks and Implicit Variational Inference

Normalizing flows

Let $\mathbf{z} \sim q_{\eta, \mathbf{x}}(\mathbf{z})$, where $q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$. Given an invertible, smooth mapping $g : \mathbb{R}^k \rightarrow \mathbb{R}^k$, the random variable $\mathbf{z}' = g(\mathbf{z})$ has a distribution

$$q(\mathbf{z}') = q_{\eta, \mathbf{x}}(\mathbf{z}) \left| \det \frac{\partial g^{-1}}{\partial \mathbf{z}'} \right| = q_{\eta, \mathbf{x}}(\mathbf{z}) \left| \det \frac{\partial g}{\partial \mathbf{z}} \right|^{-1}$$

which in general is not unimodal.

Normalizing flows

Let $\mathbf{z} \sim q_{\eta, \mathbf{x}}(\mathbf{z})$, where $q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$. Given an invertible, smooth mapping $g : \mathbb{R}^k \rightarrow \mathbb{R}^k$, the random variable $\mathbf{z}' = g(\mathbf{z})$ has a distribution

$$q(\mathbf{z}') = q_{\eta, \mathbf{x}}(\mathbf{z}) \left| \det \frac{\partial g^{-1}}{\partial \mathbf{z}'} \right| = q_{\eta, \mathbf{x}}(\mathbf{z}) \left| \det \frac{\partial g}{\partial \mathbf{z}} \right|^{-1}$$

which in general is not unimodal.

We can construct arbitrarily complex densities by composing several simple maps and applying this result.

$$\mathbf{z}_T = g_T \circ g_{T-1} \circ g_{T-2} \dots \circ g_1(\mathbf{z}), \quad \mathbf{z} \sim q_{\eta, \mathbf{x}}(\mathbf{z})$$
$$\log q_{\eta, \mathbf{x}, T}(\mathbf{z}_T) = \log q_{\eta, \mathbf{x}}(\mathbf{z}) - \sum_{t=1}^T \log \left| \det \frac{\partial g_t}{\partial \mathbf{z}_{t-1}} \right|,$$

where \mathbf{z}_t is the output of the *normalizing flow* after the t -th transformation.

Normalizing flows (II)

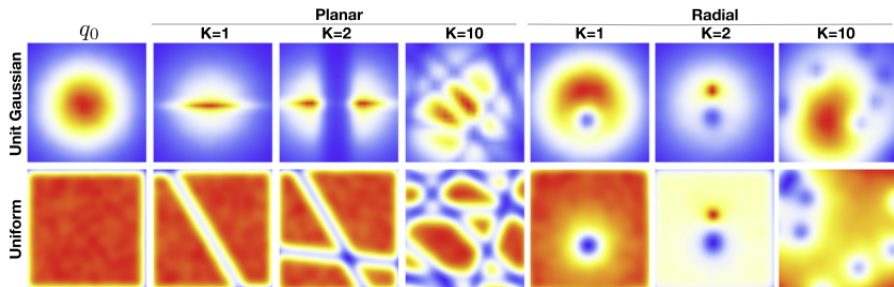


Figure 1. Effect of normalizing flow on two distributions.

Normalizing flows (III)

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}_T)p(\mathbf{z}_T)d\mathbf{z}$$

$$p(\mathbf{z}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}_T) = \mathcal{N}(\mu_{\theta}(\mathbf{z}_T), \text{diag}(\sigma_{\theta}(\mathbf{z}_T)))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\mathbf{x},T}(\mathbf{z}_T)$$

Normalizing flows (III)

$$\begin{aligned}p_{\theta}(\mathbf{x}) &= \int p_{\theta}(\mathbf{x}|\mathbf{z}_T)p(\mathbf{z}_T)d\mathbf{z} \\p(\mathbf{z}_T) &= \mathcal{N}(\mathbf{0}, \mathbf{I}_k) \\p_{\theta}(\mathbf{x}|\mathbf{z}_T) &= \mathcal{N}(\mu_{\theta}(\mathbf{z}_T), \text{diag}(\sigma_{\theta}(\mathbf{z}_T))) \\p_{\theta}(\mathbf{z}|\mathbf{x}) &\approx q_{\mathbf{x},T}(\mathbf{z}_T)\end{aligned}$$

SGD optimization is run using the following expression for the ELBO:

$$\begin{aligned}&\log p_{\theta}(\mathbf{x}) \\&\geq \mathbb{E}_{q_{\eta,\mathbf{x},T}(\mathbf{z}_T)} [\log(p_{\theta}(\mathbf{x}|\mathbf{z}_T)p(\mathbf{z}_T)) - \log q_{\eta,\mathbf{x},T}(\mathbf{z}_T)] \\&= \mathbb{E}_{q_{\eta,\mathbf{x}}(\mathbf{z})} [\log(p_{\theta}(\mathbf{x}|\mathbf{z}_T)p(\mathbf{z}_T)) - \log q_{\eta,\mathbf{x},T}(\mathbf{z}_T)] \\&= \mathbb{E}_{q_{\eta,\mathbf{x}}(\mathbf{z})} [\log(p_{\theta}(\mathbf{x}|\mathbf{z}_T)p(\mathbf{z}_T)) - \log q_{\eta,\mathbf{x}}(\mathbf{z})] + \mathbb{E}_{q_{\eta,\mathbf{x}}(\mathbf{z})} \left[\sum_{t=1}^T \log \left| \det \frac{\partial \mathbf{g}_t}{\partial \mathbf{z}_{t-1}} \right| \right]\end{aligned}$$

Index

1 Variational Autoencoders (VAEs)

- VAE Generative Model
- VAE Objective Function
- VAE Stochastic Optimization

2 Some recent applications of VAEs

3 Beyond the standard VAE generative model

- Normalizing flows
- Adversarial networks and Implicit Variational Inference

Revisiting the VAE generative model

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta,\mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \text{KL} \left(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z}) \right) \right)$$

Some critical assumptions (among many others)

- Unimodal posterior approximation.
- $p_{\theta}(\mathbf{x}|\mathbf{z})$ is known.
- I can find a valid reconstruction $p_{\theta}(\mathbf{x}|\mathbf{z})$ model.
- Prior too simple? It does not enforce interpretability in the latent space.
- ...

Implicit posterior distribution

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\boldsymbol{\eta}, \mathbf{x}}(\mathbf{z}) = f_{\boldsymbol{\eta}}(\mathbf{x}, \epsilon) \rightarrow \text{We can only sample from it}$$

$$\begin{aligned}\max_{\boldsymbol{\theta}, \boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) &= \frac{1}{N} \max_{\boldsymbol{\theta}, \boldsymbol{\eta}} \left(\sum_{i=1}^N \mathbb{E}_{q_{\boldsymbol{\eta}, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \text{KL} \left(q_{\boldsymbol{\eta}, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z}) \right) \right) \\ &= \frac{1}{N} \max_{\boldsymbol{\theta}, \boldsymbol{\eta}} \left(\sum_{i=1}^N \mathbb{E}_{q_{\boldsymbol{\eta}, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log(p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})p(\mathbf{z})) - \log(q_{\boldsymbol{\eta}, \mathbf{x}^{(i)}}(\mathbf{z})) \right] \right)\end{aligned}$$

$q_{\boldsymbol{\eta}, \mathbf{x}^{(i)}}(\mathbf{z})$ cannot appear explicitly in the objective function!

Discriminative Classifiers for log-likelihood ratio estimation

Consider the following two joint distributions:

$$\begin{cases} q_{\eta, \mathbf{x}}(\mathbf{z})p^*(\mathbf{x}) \\ p(\mathbf{z})p^*(\mathbf{x}) \end{cases} \rightarrow p^*(\mathbf{x}) \text{ is the real distribution of the data.}$$

Discriminative Classifiers for log-likelihood ratio estimation

Consider the following two joint distributions:

$$\begin{cases} q_{\eta, \mathbf{x}}(\mathbf{z})p^*(\mathbf{x}) \\ p(\mathbf{z})p^*(\mathbf{x}) \end{cases} \rightarrow p^*(\mathbf{x}) \text{ is the real distribution of the data.}$$

Let $T(\mathbf{z}, \mathbf{x})$ be a classifier network that is trained to discriminate between samples coming from either $q_{\eta, \mathbf{x}}(\mathbf{z})p^*(\mathbf{x})$ (+1 class), or $p(\mathbf{z})p^*(\mathbf{x})$ (0 class). Given a $(\mathbf{z}', \mathbf{x}')$ sample

$$\mathbb{P}((\mathbf{z}', \mathbf{x}') \text{ is drawn from } p(\mathbf{z})p^*(\mathbf{x})) = \frac{1}{1 + e^{-T(\mathbf{z}', \mathbf{x}')}} \triangleq \text{sigm}(T(\mathbf{z}', \mathbf{x}'))$$

Discriminative Classifiers for log-likelihood ratio estimation (II)

The optimal discriminator $T^*(\mathbf{z}, \mathbf{x})$ according to the cross entropy loss function

$$\arg \max_{T(\mathbf{z}, \mathbf{x})} \mathbb{E}_{q_{\eta, \mathbf{x}}(\mathbf{z}) p^*(\mathbf{x})} [\log \text{sigm}(T(\mathbf{z}, \mathbf{x}))] + \mathbb{E}_{p(\mathbf{z}) p^*(\mathbf{x})} [\log(1 - \text{sigm}(T(\mathbf{z}, \mathbf{x})))]$$

is

$$T^*(\mathbf{z}, \mathbf{x}) = \log \frac{q_{\eta, \mathbf{x}}(\mathbf{z})}{p(\mathbf{z})}$$

Discriminative Classifiers for log-likelihood ratio estimation (II)

The optimal discriminator $T^*(\mathbf{z}, \mathbf{x})$ according to the cross entropy loss function

$$\arg \max_{T(\mathbf{z}, \mathbf{x})} \mathbb{E}_{q_{\eta, \mathbf{x}}(\mathbf{z})p^*(\mathbf{x})} [\log \text{sigm}(T(\mathbf{z}, \mathbf{x}))] + \mathbb{E}_{p(\mathbf{z})p^*(\mathbf{x})} [\log(1 - \text{sigm}(T(\mathbf{z}, \mathbf{x})))]$$

is

$$T^*(\mathbf{z}, \mathbf{x}) = \log \frac{q_{\eta, \mathbf{x}}(\mathbf{z})}{p(\mathbf{z})}$$

In practice, we approximate $T^*(\mathbf{z}, \mathbf{x})$ using a deep NN with parameters ψ that is trained using **samples** of both distributions:

$$\psi^* = \arg \max_{\psi} \sum_{i=1}^N \frac{1}{2} \log \text{sigm}(T(\mathbf{x}^{(i)}, \mathbf{z}^{(i)})) + \frac{1}{2} \log(1 - \text{sigm}(T(\mathbf{x}^{(i)}, \tilde{\mathbf{z}}^{(i)})))$$

where $(\mathbf{x}^{(i)}, \mathbf{z}^{(i)})$, $i = 1, \dots, N$ are i.i.d. samples from $q_{\eta, \mathbf{x}}(\mathbf{z})p^*(\mathbf{x})$ and $(\mathbf{x}^{(i)}, \mathbf{z}^{(i)})$, $i = 1, \dots, N$ are i.i.d. samples from $p(\mathbf{z})p^*(\mathbf{x})$.

The adversarial VAE with implicit posterior approximation

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta, \mathbf{x}}(\mathbf{z}) = f_{\eta}(\mathbf{x}, \epsilon) \rightarrow \text{We can only sample from it}$$

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log(p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})) + \log\left(\frac{p(\mathbf{z})}{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})}\right) \right] \right)$$

L. Mescheder , S. Nowozin, and A. Geiger, Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks.

F. Huszár, Variational Inference using Implicit Distributions.

The adversarial VAE with implicit posterior approximation

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta, \mathbf{x}}(\mathbf{z}) = f_{\eta}(\mathbf{x}, \epsilon) \rightarrow \text{We can only sample from it}$$

$$\begin{aligned}\max_{\theta, \eta} \mathcal{L}(\theta, \eta) &= \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log(p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})) + \log\left(\frac{p(\mathbf{z})}{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})}\right) \right] \right) \\ &= \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log(p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})) - T_{\psi^*}(\mathbf{z}, \mathbf{x}^{(i)}) \right] \right)\end{aligned}$$

L. Mescheder, S. Nowozin, and A. Geiger, Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks.

F. Huszár, Variational Inference using Implicit Distributions.

The adversarial VAE with implicit posterior approximation



(a) Training data



(b) Random samples

The adversarial VAE with implicit likelihood distribution

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = f_{\theta}(\mathbf{x}, \epsilon) \rightarrow \text{We can only sample from it}$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log(p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})) + \log\left(\frac{p(\mathbf{z})}{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})}\right) \right] \right)$$

The adversarial VAE with implicit likelihood distribution

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = f_{\theta}(\mathbf{x}, \epsilon) \rightarrow \text{We can only sample from it}$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

$$\begin{aligned}\max_{\theta, \eta} \mathcal{L}(\theta, \eta) &= \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log(p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})) + \log\left(\frac{p(\mathbf{z})}{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})}\right) \right] \right) \\&= \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log\left(\frac{p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})p(\mathbf{z})}{p^*(\mathbf{x})q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})}\right) \right] \right) \\&= \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[-T_{\psi^*}(\mathbf{z}, \mathbf{x}^{(i)}) \right] \right)\end{aligned}$$

ALGORITHM	IMPLICIT			VI
	$p_{\theta}(z)$	$p_{\theta}(x z)$	$q_{\psi}(z x)$	
VAE (KINGMA & WELLING, 2014)				✓
NF (REZENDE & MOHAMED, 2015)				
PC-ADV, ALGORITHM 1				
AFFGAN [†] (SØNDERBY ET AL., 2017)	I		✓	✓
AVB (MESCHEDER ET AL., 2017)				
OPVI (RANGANATH ET AL., 2016)	I		✓	✓
PC-DEN, ALGORITHM 3	I		✓	✓
JC-ADV, ALGORITHM 2	I	I	✓	✓
JC-DEN	I	I	✓	✓
JC-ADV-RMD [‡]	✓	✓	✓	✓
AAE (MAKHZANI ET AL., 2016)	I		✓	
DEEPSIM (DOSOVITSKIY & BROX, 2016)	I		✓	
ALI (DUMOULIN ET AL., 2017)	✓	✓	✓	
BiGAN (DONAHUE ET AL., 2017)				