

Variational Auto-Encoders

A probabilistic introduction



Alberto Lumbreras

a.lumbreras@criteo.com

April 2020

1. Two schools of thought

2. Maximum likelihood estimation

- Basic MLE
- MLE with Latent Variables

3. Variational Autoencoders

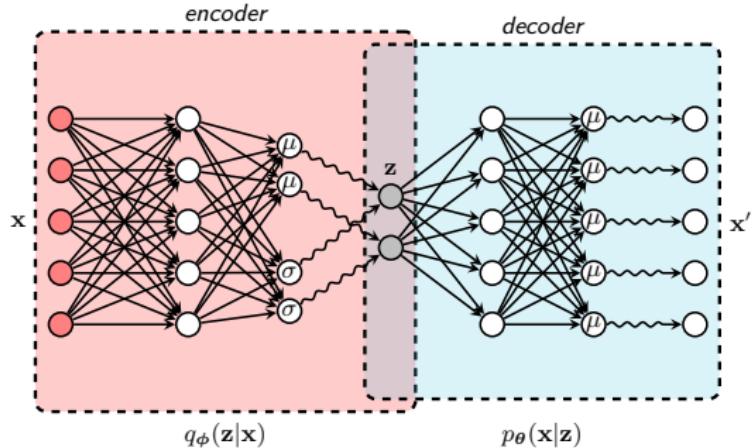
- You cannot always push gradients inside expectations
- The old and the new hack
- Scaling with amortization
- The trinity of the Evidence Lower Bound

4. VAE algorithms

- Gaussian VAE
- Pseudo-code to train
- Pseudo-code to generate synthetic samples
- Pseudo-code to generate latent representations

5. FAQ

Variational Autoencoder overview



$$\begin{aligned}
 \mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \\
 &\approx \sum_j \log p_{\theta}(\mathbf{x}|\mathbf{z}^{(j)}) - \log q_{\phi}(\mathbf{z}^{(j)}|\mathbf{x}) - \log p_{\theta}(\mathbf{z}^{(j)}), \quad \mathbf{z}^{(j)} \sim q_{\phi}(\mathbf{z}|\mathbf{x}) \\
 &\approx \sum_j \log p_{\theta}(\mathbf{x}|\mathbf{z}^{(j)}) - \log q_{\phi}(\mathbf{z}^{(j)}|\mathbf{x}) - \log p_{\theta}(\mathbf{z}^{(j)}), \quad \mathbf{z}^{(j)} = g_{\phi}(\boldsymbol{\epsilon}^{(j)}, \mathbf{x}), \\
 &\quad \boldsymbol{\epsilon}^{(j)} \sim p(\boldsymbol{\epsilon})
 \end{aligned}$$

Two schools of thought



- Connectionist school:

- An auto-encoder that uses Variational Inference.
- Scalable method for data generation.
- Competes with Generative Adversarial Networks (GANs).
- Mathematical justification at the end.



- Probabilistic school:

- A Variational Inference that looks like an auto-encoder.
- Scalable method for density estimation.
- Competes with of Mean-Field Variational Bayes (MF-VB).
- Mathematical justification from the beginning.



1. Two schools of thought

2. Maximum likelihood estimation

- Basic MLE
- MLE with Latent Variables

3. Variational Autoencoders

- You cannot always push gradients inside expectations
- The old and the new hack
- Scaling with amortization
- The trinity of the Evidence Lower Bound

4. VAE algorithms

- Gaussian VAE
- Pseudo-code to train
- Pseudo-code to generate synthetic samples
- Pseudo-code to generate latent representations

5. FAQ

1. Two schools of thought

2. Maximum likelihood estimation

- Basic MLE
- MLE with Latent Variables

3. Variational Autoencoders

4. VAE algorithms

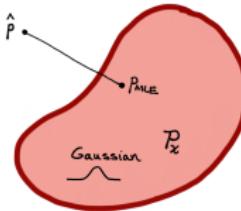
5. FAQ

Maximum Likelihood Estimation (MLE)



- We observe $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^N$, i.i.d samples from some (unknown) distribution \hat{p} ,
- $$\mathbf{x}_n \sim \hat{p}(\mathbf{x})$$

- We want to find the distribution $p_{\theta} \in \mathcal{P}_x$ that is closest to \hat{p} .



- The Maximum Likelihood Estimator minimises the Kullback-Leibler divergence:

$$\begin{aligned}
 D_{KL}(\hat{p}(\mathbf{x}) \| p_{\theta}(\mathbf{x})) &= \mathbb{E}_{\hat{p}(\mathbf{x})} \left[\log \frac{\hat{p}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right] \\
 &= -\mathbb{E}_{\hat{p}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})] + \text{const.} \\
 &\approx -\sum_{n=1}^N \log p_{\theta}(\mathbf{x}_n) + \text{const.}
 \end{aligned}$$



1. Two schools of thought

2. Maximum likelihood estimation

- Basic MLE
- MLE with Latent Variables

3. Variational Autoencoders

4. VAE algorithms

5. FAQ

Maximum Likelihood Estimation with Latent Variables

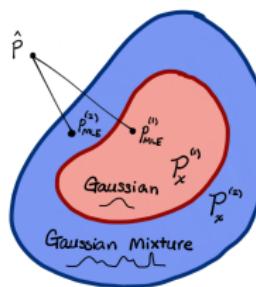


- We can assume each observation depends on a latent (non-observed) variable \mathbf{z} :

$$\mathbf{z}_n \sim \hat{p}(\mathbf{z})$$

$$\mathbf{x}_n \sim \hat{p}(\mathbf{x}|\mathbf{z}_n)$$

- This increases expressiveness,



- But complicates Maximum Likelihood Estimation:

$$\arg \max_{\theta} \sum_{n=1}^N \log p_{\theta}(\mathbf{x}_n) = \arg \max_{\theta} \sum_{n=1}^N \log \int p_{\theta}(\mathbf{x}_n, \mathbf{z}_n) d\mathbf{z}_n$$

The intractable integral and the Evidence Lower Bound



- ① The **integral** is often intractable. We can re-express it as:

$$\log p_{\theta}(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \log \int \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z})} q_{\phi}(\mathbf{z}) d\mathbf{z} = \log \mathbb{E}_{q_{\phi}(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z})} \right]$$

where $q_{\phi}(\mathbf{z})$ can be **any** density¹.

- ② Use Jensen ($\log \mathbb{E}_{q(z)}[f(z)] \geq \mathbb{E}_{q(z)}[\log f(z)]$) to find a **tractable lower bound**:

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z})} \right] = \mathcal{L}_{\theta, \phi}(\mathbf{x}) \quad (1)$$



Posterior and density estimation

$$\log p_{\theta}(\mathbf{x}) = \mathcal{L}_{\theta, \phi}(\mathbf{x}) \iff q_{\phi}(\mathbf{z}) = p_{\theta}(\mathbf{z}|\mathbf{x})$$

¹For those familiar with Importance Sampling: it is the same expression.

Maximizing the ELBO

Strategies for posteriors in close form, samples or approximations.



- Expectation-Maximization (EM) ([Dempster et al. 1977](#)):
(posterior: analytically available)

$$\mathcal{L}_{\theta}(\mathbf{X}) = \mathbb{E}_{p_{\theta}(\mathbf{Z}|\mathbf{X})}[\log p_{\theta}(\mathbf{X}, \mathbf{Z}) - \log p_{\theta}(\mathbf{Z}|\mathbf{X})]$$



- Monte-Carlo EM (MC-EM) ([Wei and Tanner 1990](#)):
(posterior: analytically not available)

$$\mathcal{L}_{\theta}(\mathbf{X}) = \frac{1}{J} \sum_j \{\log p_{\theta}(\mathbf{X}, \mathbf{Z}^{(j)}) - \log p_{\theta}(\mathbf{Z}^{(j)}|\mathbf{X})\}, \quad \mathbf{Z}^{(j)} \sim p_{\theta}(\mathbf{Z}|\mathbf{X})$$



- Mean-Field Variational Bayes EM (MFVB-EM) ([Jordan et al. 1999](#)):
(posterior: approximation analytically available)

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z})}[\log p_{\theta}(\mathbf{X}, \mathbf{Z}) - \log q_{\phi}(\mathbf{Z})], \quad q_{\phi}(\mathbf{Z}) = \prod_n q_{\phi}^{\{n\}}(\mathbf{z}_n)$$



None of the above methods can be used for complicated models

1. Two schools of thought
2. Maximum likelihood estimation
 - Basic MLE
 - MLE with Latent Variables
3. **Variational Autoencoders**
 - You cannot always push gradients inside expectations
 - The old and the new hack
 - Scaling with amortization
 - The trinity of the Evidence Lower Bound
4. VAE algorithms
 - Gaussian VAE
 - Pseudo-code to train
 - Pseudo-code to generate synthetic samples
 - Pseudo-code to generate latent representations
5. FAQ

1. Two schools of thought

2. Maximum likelihood estimation

3. Variational Autoencoders

- You cannot always push gradients inside expectations
- The old and the new hack
- Scaling with amortization
- The trinity of the Evidence Lower Bound

4. VAE algorithms

5. FAQ

Maximizing the ELBO without assumptions

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z})} [\log_{\theta} p(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})]$$



- The posterior is too complex to make simplistic independence assumptions (MFVB-EM)



A very generic strategy:

- ① Push the gradient inside the expectation (careful: *here be dragons*)².
- ② And **then** we approximate the expectation by Monte-Carlo sampling.

²See Mohamed et al. 2019 for an excellent review on this problem.

Pushing gradients inside expectations

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z})} [\log_{\theta} p(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})]$$



- The gradient w.r.t. θ can be pushed-in:



$$\nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\theta}(\mathbf{z})] = \nabla_{\theta} \int f_{\theta}(\mathbf{z}) q_{\phi}(\mathbf{z}) d\mathbf{z} \quad (2)$$

$$= \int \{\nabla_{\theta} f_{\theta}(\mathbf{z})\} q_{\phi}(\mathbf{z}) d\mathbf{z} \quad (3)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z})} [\nabla_{\theta} f_{\theta}(\mathbf{z})] \approx \frac{1}{J} \sum_j \nabla_{\theta} f_{\theta}(\mathbf{z}^{(j)}) \quad (4)$$

- The gradient w.r.t. the posterior parameters ϕ cannot be pushed-in:³



$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\theta}(\mathbf{z})] = \nabla_{\phi} \int f_{\theta}(\mathbf{z}) q_{\phi}(\mathbf{z}) d\mathbf{z} \quad (5)$$

$$= \int f_{\theta}(\mathbf{z}) \nabla_{\phi} q_{\phi}(\mathbf{z}) d\mathbf{z} \quad (6)$$

$$\neq \mathbb{E}_{q_{\phi}(\mathbf{z})} [\nabla_{\phi} f_{\theta}(\mathbf{z})] \approx \frac{1}{J} \sum_j \nabla_{\phi} f_{\theta}(\mathbf{z}^{(j)}) \quad (7)$$

³This problem did not appear in the era before Variational Inference, since we only had θ .

1. Two schools of thought

2. Maximum likelihood estimation

3. Variational Autoencoders

- You cannot always push gradients inside expectations
- **The old and the new hack**
- Scaling with amortization
- The trinity of the Evidence Lower Bound

4. VAE algorithms

5. FAQ

The old hack

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z})} [\log_{\theta} p(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})]$$



- The score function:

$$\nabla_{\theta} \log p_{\theta}(x) = \frac{\nabla_{\theta} p_{\theta}(x)}{p_{\theta}(x)} \implies \nabla_{\theta} p_{\theta}(x) = p_{\theta}(x) \nabla_{\theta} \log p_{\theta}(x) \quad (8)$$

- The gradient w.r.t. the posterior parameters ϕ ~~cannot be pushed in~~:

$$\begin{aligned} \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\theta}(\mathbf{z})] &= \nabla_{\phi} \int f_{\theta}(\mathbf{z}) q_{\phi}(\mathbf{z}) d\mathbf{z} \\ &= \int f_{\theta}(\mathbf{z}) \nabla_{\phi} q_{\phi}(\mathbf{z}) d\mathbf{z} \\ &= \int f_{\theta}(\mathbf{z}) \mathbf{q}_{\phi}(\mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z})} [f_{\theta}(\mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z})] \\ &\approx \frac{1}{J} \sum_j f_{\theta}(\mathbf{z}^{(j)}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}^{(j)}) \end{aligned}$$



- Needs variance reduction (Paisley et al. 2012).

The new hack: Stochastic Gradient Variational Bayes estimator (SGVB)



- The Law Of the Unconscious Statistician (LOTUS):

Given a random variable ϵ with a probability distribution $p(\epsilon)$ and a transformation $z = g(\epsilon)$:⁴

$$\mathbb{E}_{q_\phi(z)}[f(z)] = \mathbb{E}_{p(\epsilon)}[f(g_\phi(\epsilon))] \quad (9)$$

- The gradient w.r.t. the posterior parameters ϕ ~~cannot be pushed in~~:

$$\begin{aligned} \nabla_\phi \mathbb{E}_{q_\phi(z)}[f_\theta(z)] &= \nabla_\phi \mathbb{E}_{p(\epsilon)}[f_\theta(g_\phi(\epsilon))] \\ &= \mathbb{E}_{p(\epsilon)}[\nabla_\phi f_\theta(g_\phi(\epsilon))] \\ &\approx \frac{1}{J} \sum_j \nabla_\phi f_\theta(g_\phi(\epsilon^{(j)})) \end{aligned}$$



⁴For those familiar with Normalizing Flows: $g(\epsilon)$ can be a composition $g_k(g_{k-1}(\dots g_1(\epsilon)))$

Stochastic Gradient Variational Bayes estimator (SGVB)

Using LOTUS and amortized inference



- Stochastic Gradient Variational Bayes estimator:
(Kingma and Welling 2014)

$$\epsilon^{(j)} \sim p(\epsilon)$$

$$\mathbf{z}^{(j)} = g_\phi((\epsilon^{(j)}, \mathbf{x}^{(j)}))$$

$$\nabla_\phi \mathcal{L}_{\theta, \phi}(\mathbf{x}) \approx \frac{1}{J} \sum_j \nabla_\phi \log p_\theta(\mathbf{x}, \mathbf{z}^{(j)}) - \nabla_\phi \log q_\phi(\mathbf{z}^{(j)} | \mathbf{x})$$

Samples are drawn from a distribution that is not affected by the gradient

1. Two schools of thought

2. Maximum likelihood estimation

3. Variational Autoencoders

- You cannot always push gradients inside expectations
- The old and the new hack
- **Scaling with amortization**
- The trinity of the Evidence Lower Bound

4. VAE algorithms

5. FAQ

Amortized Variational Inference

Making it scalable by reducing the number of parameters⁵



- Classic Variational Inference:

- Posterior $q(\mathbf{z}_n; \phi_n)$ depends on local parameters.
- Variational parameters: $\mathcal{O}(N)$.

- Amortized Variational Inference:

- Posterior $q(\mathbf{z}_n; f_\phi(\mathbf{x}_n))$ depends on shared parameters (neural net), .
- Variational parameters: $\mathcal{O}(1)$.

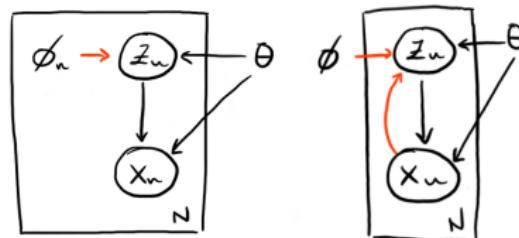


Figure: Non-amortized (left) and amortized variational models (right). Red lines represent the dependencies of the variational distribution of \mathbf{z}_n

⁵See Zhang et al. 2019 for more context.

1. Two schools of thought

2. Maximum likelihood estimation

3. Variational Autoencoders

- You cannot always push gradients inside expectations
- The old and the new hack
- Scaling with amortization
- **The trinity of the Evidence Lower Bound**

4. VAE algorithms

5. FAQ

The trinity of the ELBO



- The classic:
(expected joint log-likelihood plus entropy)

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log q_\phi(\mathbf{z}|\mathbf{x})] \quad (10)$$

- The variational:
(evidence minus quality of the posterior approximation)

$$p(\mathbf{x}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \quad (11)$$

- The auto-encoder:
(quality of reconstruction minus regularization term)

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (12)$$

The three are the same.



Why is it called Variational Autoencoder

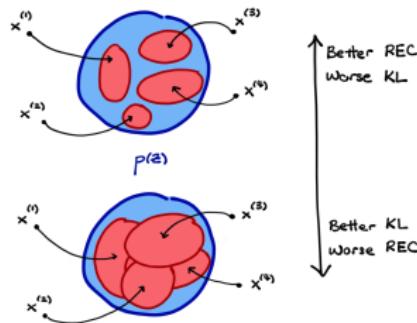


- The third form of the ELBO trinity is:

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

where we identified a **reconstruction** and a **regularizer** term.

- This view makes explicit the trade-off between fitting the data (reconstruction) and fitting the prior (regularization):⁶



⁶The nice drawings are taken from <http://ruishu.io/2018/03/14/vae/>

1. Two schools of thought
2. Maximum likelihood estimation
 - Basic MLE
 - MLE with Latent Variables
3. Variational Autoencoders
 - You cannot always push gradients inside expectations
 - The old and the new hack
 - Scaling with amortization
 - The trinity of the Evidence Lower Bound
4. VAE algorithms
 - Gaussian VAE
 - Pseudo-code to train
 - Pseudo-code to generate synthetic samples
 - Pseudo-code to generate latent representations
5. FAQ

1. Two schools of thought
2. Maximum likelihood estimation
3. Variational Autoencoders
4. **VAE algorithms**
 - Gaussian VAE
 - Pseudo-code to train
 - Pseudo-code to generate synthetic samples
 - Pseudo-code to generate latent representations
5. FAQ

Variational Auto-Encoder

Gaussian posterior, Gaussian likelihood



- For the re-parameterization trick to be used, we need:

- A family distribution for the posterior $q_\phi(z|x)$.
- A differentiable transformation $z = g_\phi(\epsilon)$.
- An distribution over the auxiliary variable $p(\epsilon)$

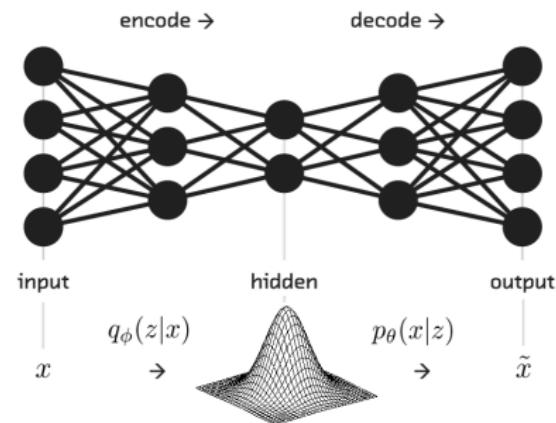
- Gaussian posterior (encoder):

- $z|x \sim \mathcal{N}(z|f_\mu(x), f_\sigma(x)\mathbf{I})$
- $z = g_\phi(x, \epsilon) = f_\mu^{(e)}(x) + f_\sigma^{(e)}(x)\mathbf{I}\epsilon$
- $\epsilon \sim \mathcal{N}(0, 1)$

- Gaussian likelihood (decoder):

- $x|z \sim \mathcal{N}(x|f_\mu(z), f_\sigma(z)\mathbf{I})$
- $x = g_\theta(z, \epsilon) = f_\mu^{(d)}(z) + f_\sigma^{(d)}(z)\mathbf{I}\epsilon$
- $\epsilon \sim \mathcal{N}(0, 1)$

where $f^{(e)}, f^{(d)}$ can be any differentiable function such as a neural network.



1. Two schools of thought
2. Maximum likelihood estimation
3. Variational Autoencoders
4. **VAE algorithms**
 - Gaussian VAE
 - **Pseudo-code to train**
 - Pseudo-code to generate synthetic samples
 - Pseudo-code to generate latent representations
5. FAQ

Training algorithm



Algorithm 1: Batch version of Auto-Encoding VB

Input: Observed data vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$

Parameters: Base distribution $p(\epsilon)$, sampling path $\mathbf{z} = g_\phi(\epsilon, \mathbf{x})$,
prior distribution $p(\mathbf{z})$, likelihood $p_\theta(\mathbf{x}|\mathbf{z})$.

Output: Model parameters θ , variational parameters ϕ

Initialize: Random initialization of θ and ϕ

repeat

 Sample $\epsilon_n \sim p(\epsilon)$

$\mathbf{z}_n = g_\phi(\epsilon, \mathbf{x}_n)$

 Compute gradients:

$$\frac{1}{N} \sum_n \nabla_{\theta, \phi} [\log p_\theta(\mathbf{x}_n | \mathbf{z}_n) - \{\log q_\phi(\mathbf{z}_n | \mathbf{x}_n) - \log p(\mathbf{z}_n)\}]$$

$(\theta, \phi) \leftarrow$ Update parameters using gradients (e.g.: SGD)

until convergence of (θ, ϕ) ;

- The more frequently used mini-batch version can be easily derived from this.
- The parameters of the encoder q_ϕ are the outputs of a neural net with inputs \mathbf{x} .
- The parameters of the decoder p_θ are the outputs of a neural net with inputs \mathbf{z} .

1. Two schools of thought
2. Maximum likelihood estimation
3. Variational Autoencoders
4. VAE algorithms
 - Gaussian VAE
 - Pseudo-code to train
 - **Pseudo-code to generate synthetic samples**
 - Pseudo-code to generate latent representations
5. FAQ

Algorithm to generate synthetic samples



- We want to sample from

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}. \quad (13)$$

Algorithm 2: VAE sample generator

Parameters: Prior distribution $p(\mathbf{z})$, likelihood $p_{\theta}(\mathbf{x}|\mathbf{z})$, decoder parameters θ

Output: Synthetic data points \mathbf{x}

Sample $\mathbf{z} \sim p(\mathbf{z})$

Sample $\mathbf{x} \sim p_{\theta}(\mathbf{x}|\mathbf{z})$

1. Two schools of thought
2. Maximum likelihood estimation
3. Variational Autoencoders
4. **VAE algorithms**
 - Gaussian VAE
 - Pseudo-code to train
 - Pseudo-code to generate synthetic samples
 - **Pseudo-code to generate latent representations**
5. FAQ

Algorithm to generate latent representation



- We want to sample from

$$p_{\theta}(z|x) \quad (14)$$

Algorithm 3: VAE code generator

Input: Observed data vector x

Parameters: Posterior distribution $q_{\phi}(z|x)$, encoder parameters ϕ

Output: Latent vector z

Sample $z \sim p_{\phi}(z|x)$

1. Two schools of thought
2. Maximum likelihood estimation
 - Basic MLE
 - MLE with Latent Variables
3. Variational Autoencoders
 - You cannot always push gradients inside expectations
 - The old and the new hack
 - Scaling with amortization
 - The trinity of the Evidence Lower Bound
4. VAE algorithms
 - Gaussian VAE
 - Pseudo-code to train
 - Pseudo-code to generate synthetic samples
 - Pseudo-code to generate latent representations
5. FAQ

Frequently Asked Questions



- **Do we loose expressiveness when we do Amortized Inference?**

Yes. This is known as the amortization error ([Cremer et al. 2018](#)). In some cases ([Sønderby et al. 2016](#)), this gap is reasonably small. In general, the more layers your neural nets have, the shorter this gap can be. Note also that this gap is the only thing that prevents the ELBO from reaching $p(x)$. In other words, without this gap, the KL divergence between the variational and the true posterior could be zero.

- **Can we use discrete likelihoods for discrete data?**

Of course. For instance, for binary data you can use a Bernoulli likelihood.

- **Can we use discrete latent variables?**

Not straightforwardly, but yes, with some hack called the Gumbel-Max trick ([Jang et al. 2017](#)).

- **What should I use as a prior $p(z)$?**

People often use a standard normal distribution. You can use whatever you want. Of course there will be better and worse priors for each problem. As there are better and worse regularization terms in any other model.

- **Can I sample from the approximate posterior $q(z|x)$?**

Yes. Just put an x in the encoder and it will output the parameters of the posterior (e.g.: the mean and variance of a Gaussian, if we choose our posterior to be Gaussian) to sample from.

References |



-  Cremer, Chris et al. (2018). "Inference Suboptimality in Variational Autoencoders". In: *Proceedings of the 35th International Conference on Machine Learning (ICML '18)*. Vol. 80, pp. 1078–1086.
-  Dempster, A. P. et al. (1977). "Maximum Likelihood from Incomplete Data via the EM Algorithm". In: *Journal of the Royal Statistical Society. Series B (Methodological)* 39.1, pp. 1–38.
-  Jang, Eric et al. (2017). "Categorical Reparametrization with Gumbel-Softmax". In: *Proceedings International Conference on Learning Representations (ICLR '17)*.
-  Jordan, Michael I. et al. (1999). "An Introduction to Variational Methods for Graphical Models". In: *Machine Learning* 37.2, pp. 183–233.
-  Kingma, Diederik P. and Max Welling (2014). "Auto-Encoding Variational Bayes". In: *2nd International Conference on Learning Representations, (ICLR '14)*.
-  Mohamed, Shakir et al. (2019). *Monte Carlo Gradient Estimation in Machine Learning*. arXiv: [1906.10652 \[stat.ML\]](https://arxiv.org/abs/1906.10652).
-  Paisley, John et al. (2012). "Variational Bayesian Inference with Stochastic Search". In: *Proceedings of the 29th International Conference on International Conference on Machine Learning (ICML '12)*, pp. 1363–1370.

References II



-  Sønderby, Casper Kaae et al. (2016). "Ladder Variational Autoencoders". In: *Advances in Neural Information Processing Systems 29*, pp. 3738–3746.
-  Wei, Greg C. G. and Martin A. Tanner (1990). "A Monte Carlo Implementation of the EM Algorithm and the Poor Man's Data Augmentation Algorithms". In: *Journal of the American Statistical Association* 85.411, pp. 699–704.
-  Zhang, C. et al. (2019). "Advances in Variational Inference". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41.08, pp. 2008–2026.

Appendix



Change of variable



- Let us consider density $q_\phi(\mathbf{z})$ and a change of variable $\mathbf{z} = g_\phi(\epsilon)$.
- The probability contained in a differential area must be invariant under change of variables.

$$p(\epsilon) d\epsilon = p(\mathbf{z}) d\mathbf{z} \quad (15)$$

$$p(\epsilon) d\epsilon = p(g(\epsilon)) dg(\epsilon) \quad (16)$$

$$p(\epsilon) = p(g(\epsilon)) \left| \frac{dg(\epsilon)}{d\epsilon} \right| \quad (17)$$

- Application to expectations (LOTUS):

$$\mathbb{E}_{q_\phi(\mathbf{z})}[f(\mathbf{z})] = \int f(\mathbf{z}) q_\phi(\mathbf{z}) d\mathbf{z} = \int f(g_\phi(\epsilon)) p(\epsilon) d\epsilon = \mathbb{E}_{p(\epsilon)}[f(g_\phi(\epsilon))] \quad (18)$$

Pushing gradients inside expectations

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z})} [\log_{\theta} p(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})]$$



- The gradient w.r.t. θ can be pushed-in:⁷:

$$\nabla_{\theta} \mathcal{L}_{\theta, \phi}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] = \nabla_{\theta} \int \log p_{\theta}(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z}) d\mathbf{z} \quad (19)$$

$$= \int \{\nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z})\} q_{\phi}(\mathbf{z}) d\mathbf{z} = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z})] \quad (20)$$

- The gradient w.r.t. the posterior parameters ϕ cannot be pushed-in:

$$\nabla_{\phi} \mathcal{L}_{\theta, \phi}(\mathbf{x}) = \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})] \quad (21)$$

$$= \nabla_{\phi} \int [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})] q_{\phi}(\mathbf{z}) d\mathbf{z} \quad (22)$$

$$= \int \nabla_{\phi} \{[\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})] q_{\phi}(\mathbf{z})\} d\mathbf{z} \quad (23)$$

$$\neq \mathbb{E}_{q_{\phi}(\mathbf{z})} [\nabla_{\phi} \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})] \quad (24)$$

⁷For those familiar with Expectation Maximization: this gradient appears in the M-step.
This push-in is used in MC-EM. No need to push-in if we can solve the expectation analytically.

Useful identities



- Score function:

$$\nabla_{\theta} \log p_{\theta}(x) = \frac{\nabla_{\theta} p_{\theta}(x)}{p_{\theta}(x)} \quad (25)$$

- The Law Of the Unconscious Statistician (LOTUS):

Given a random variable ϵ with a probability distribution $p(\epsilon)$ and a transformation $z = g(\epsilon)$:⁸

$$\mathbb{E}_{p(z)}[z] = \mathbb{E}_{p(\epsilon)}[g(\epsilon)] \quad (26)$$

which means that we can compute the expectation of a function of a random variable z without knowing its distribution, if we know its corresponding sampling path and base distribution.

⁸For those familiar with Normalizing Flows: $g(\epsilon)$ can be a composition $g_k(g_{k-1}(\dots g_1(\epsilon)))$

Naive Monte-Carlo gradient estimator

Using the Score Function



- Naïve Monte-Carlo gradient estimator (Paisley et al. 2012):

Using the score function, we obtain:

$$\begin{aligned}\nabla_{\phi} \mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z})} [\{\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})\} \nabla_{\phi} \log q_{\phi}(\mathbf{z})] \\ &\approx \frac{1}{J} \sum_j \{\log p_{\theta}(\mathbf{x}, \mathbf{z}^{(j)}) - \log q_{\phi}(\mathbf{z}^{(j)})\} \nabla_{\phi} \log q_{\phi}(\mathbf{z}^{(j)})\end{aligned}$$

where $\mathbf{z}^{(j)} \sim q_{\phi}(\mathbf{z})$ (27)

This estimator suffers from a high variance