DeepMind

Unbiased Gradient Estimation for Variational **Auto-Encoders using Coupled Markov Chains**

Francisco J. R. Ruiz, Michalis K. Titsias, Taylan Cemgil, Arnaud Doucet

6

24/11/2020

DeepMind

Motivation and Goal

Private & Confidential

Deep Generative Models

Variational Auto-Encoder (VAE)







Low resolution

High resolution

[Anonymous, ICLR 2021]

Energy-Based Model (EBM)







[OpenAl]



Fitting Deep Generative Models

- Stochastic gradient descent (SGD) is a **powerful tool** in Machine Learning
- SGD obtains and follows **unbiased estimates** of the log-likelihood gradient
- For some deep generative models, unbiased gradient estimates are **not directly available**





Why Fitting These Models is Hard



(but it can be written in terms of an expectation)

Summary of Contributions

- An algorithm to obtain **unbiased gradient estimates** for VAEs
 - **Avoid approximation gap** of previous approaches (based on bounds)
 - Lead to VAEs with **better predictive performance**
 - Applicable to other deep generative models beyond VAEs
 - Main **limitation**: increased computational complexity



Technical Tools

- Two main ideas to develop the unbiased gradient estimates:
 - Augmented latent space (similarly to IWAE)
 - MCMC couplings based on importance sampling





DeepMind

Preliminaries

Review: VAE / IWAE

• The VAE log-likelihood



• Optimize the ELBO (a lower bound)

$$\mathcal{L}_{ ext{ELBO}}(heta, \phi) = \mathbb{E}_{q_{\phi}(z \mid x)} \left[\log w_{ heta, \phi}(z)
ight]$$



auxiliary random variables

• Or optimize the IWAE (a tighter lower bound)

$$\mathcal{L}_{\text{IWAE}}(\theta, \phi) = \mathbb{E}_{q_{\phi}(z_{1:K} \mid x)} \left[\log \left(\frac{1}{K} \sum_{k=1}^{K} w_{\theta, \phi}(z_k) \right) \right]$$

Private & Confidential





Background: The IWAE as an Augmented Space



Private & Confidential

The Roadmap to Unbiased Estimation

- Both the ELBO and IWAE bounds are **biased** approximators of gradient of the log-likelihood
- If we could **sample** from the posterior, we could easily form an unbiased estimator

expectation /





gradient

• MCMC couplings provide unbiased estimators by design without the need to obtain exact samples from the posterior

 $\nabla_{\theta} \mathcal{L}(\theta) = \mathbb{E}_{p_{\theta}(z \mid x)} \left[\nabla_{\theta} \log p_{\theta}(x, z) \right]$



MCMC Couplings

• Consider estimating an expectation of the form





• Typical MCMC approach: sample from a kernel $\mathcal{K}(\cdot \mid u)$ that targets $\pi(u)$



Private & Confidential



MCMC Couplings

• Consider estimating an expectation of the form

$$H = \mathbb{E}_{\pi(u)}[h(u)]$$

generic r.v.



- Typical MCMC approach: sample from a kernel $\mathcal{K}(\cdot \mid u)$ that targets $\pi(u)$
- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \overline{u})$

MCMC Couplings

Consider estimating an expectation of the form

$$H = \mathbb{E}_{\pi(u)}[h(u)]$$

generic r.v.



- Typical MCMC approach: sample from a kernel $\mathcal{K}(\cdot \mid u)$ that targets $\pi(u)$
- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \bar{u})$





- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - Marginally evolving according to $\mathcal{K}(\cdot \mid u)$
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \overline{u})$





- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - Marginally evolving according to $\mathcal{K}(\cdot \mid u)$
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \overline{u})$
 - Initialize the first chain $u^{(1)} \sim \mathcal{K}(u \mid u^{(0)})$





- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - Marginally evolving according to $\mathcal{K}(\cdot \mid u)$
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \overline{u})$
 - Initialize the first chain $u^{(1)} \sim \mathcal{K}(u \mid u^{(0)})$
 - Then sample both chains from the joint kernel $u^{(t+1)}, \bar{u}^{(t)} \sim \mathcal{K}_C(u, \bar{u} \mid u^{(t)}, \bar{u}^{(t-1)})$





- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - Marginally evolving according to $\mathcal{K}(\cdot \mid u)$
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \overline{u})$
 - Initialize the first chain $u^{(1)} \sim \mathcal{K}(u \mid u^{(0)})$
 - Then sample both chains from the joint kernel $u^{(t+1)}, \bar{u}^{(t)} \sim \mathcal{K}_C(u, \bar{u} \mid u^{(t)}, \bar{u}^{(t-1)})$





- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - Marginally evolving according to $\mathcal{K}(\cdot \mid u)$
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \overline{u})$
 - Initialize the first chain $u^{(1)} \sim \mathcal{K}(u \mid u^{(0)})$
 - Then sample both chains from the joint kernel $u^{(t+1)}, \bar{u}^{(t)} \sim \mathcal{K}_C(u, \bar{u} \mid u^{(t)}, \bar{u}^{(t-1)})$





- Coupling MCMC: use **two** MCMC chains with the same stationary distribution that are *coupled*
 - Marginally evolving according to $\mathcal{K}(\cdot \mid u)$
 - There is a joint MCMC kernel $\mathcal{K}_C(\cdot, \cdot \mid u, \overline{u})$
 - Initialize the first chain $u^{(1)} \sim \mathcal{K}(u \mid u^{(0)})$
 - Then sample both chains from the joint kernel $u^{(t+1)}, \bar{u}^{(t)} \sim \mathcal{K}_C(u, \bar{u} \mid u^{(t)}, \bar{u}^{(t-1)})$
- Define the meeting time $\tau = \inf\{t \ge 1 : u^{(t)} = \overline{u}^{(t-1)}\}$
 - We design the kernel such that: (i) the two chains meet each other (τ is random but finite), and (ii) they remain equal to each other afterwards





• Use the samples from both chains to form an estimator

$$H \approx h(u^{(t_0)}) + \sum_{t=t_0+1}^{\tau-1} \left(h(u^{(t)}) - h(\bar{u}^{(t-1)}) \right)$$



$$\tau = \inf\{t \ge 1 : u^{(t)} = \bar{u}^{(t-1)}\}$$



• Use the samples from both chains to form an estimator

$$H \approx h(u^{(t_0)}) + \sum_{t=t_0+1}^{\tau-1} \left(h(u^{(t)}) - h(\bar{u}^{(t-1)}) \right)$$

$$\begin{array}{c} u^{(1)} & & u^{(2)} & & u^{(3)} \\ \hline u^{(1)} & & & \overline{u}^{(2)} & & \overline{u}^{(3)} \end{array}$$

 $\tau = \inf\{t \ge 1 : u^{(t)} = \bar{u}^{(t-1)}\}$

Proof

$$H = \mathbb{E}_{\pi(u)}[h(u)] = \mathbb{E}[h(u^{(t_0)})] + \sum_{t=t_0+1}^{\infty} \left(\mathbb{E}[h(u^{(t)})] - \mathbb{E}[h(u^{(t-1)})] \right)$$
 telescoping sum

• Use the samples from both chains to form an estimator

$$H \approx h(u^{(t_0)}) + \sum_{t=t_0+1}^{\tau-1} \left(h(u^{(t)}) - h(\bar{u}^{(t-1)}) \right)$$



 $\tau = \inf\{t \ge 1 : u^{(t)} = \bar{u}^{(t-1)}\}$

H

$$\begin{split} &= \mathbb{E}_{\pi(u)}[h(u)] = \mathbb{E}[h(u^{(t_0)})] + \sum_{t=t_0+1}^{\infty} \left(\mathbb{E}[h(u^{(t)})] - \mathbb{E}[h(u^{(t-1)})] \right) & \quad telescoping \ \text{sum} \\ &= \mathbb{E}[h(u^{(t_0)})] + \sum_{t=t_0+1}^{\infty} \left(\mathbb{E}[h(u^{(t)})] - \mathbb{E}[h(\overline{u}^{(t-1)})] \right) & \quad \text{same marginals} \end{split}$$



• Use the samples from both chains to form an estimator

$$H \approx h(\boldsymbol{u}^{(t_0)}) + \sum_{t=t_0+1}^{\tau-1} \left(h(\boldsymbol{u}^{(t)}) - h(\bar{\boldsymbol{u}}^{(t-1)}) \right)$$



 $\tau = \inf\{t \ge 1 : u^{(t)} = \bar{u}^{(t-1)}\}$

$$\begin{split} H &= \mathbb{E}_{\pi(u)}[h(u)] = \mathbb{E}[h(u^{(t_0)})] + \sum_{t=t_0+1}^{\infty} \left(\mathbb{E}[h(u^{(t)})] - \mathbb{E}[h(u^{(t-1)})] \right) & \text{telescoping sum} \\ &= \mathbb{E}[h(u^{(t_0)})] + \sum_{t=t_0+1}^{\infty} \left(\mathbb{E}[h(u^{(t)})] - \mathbb{E}[h(\bar{u}^{(t-1)})] \right) & \text{same marginals} \\ &= \mathbb{E}\left[h(u^{(t_0)}) + \sum_{t=t_0+1}^{\infty} \left(h(u^{(t)}) - h(\bar{u}^{(t-1)}) \right) \right] & \text{swap expectation and limit} \end{split}$$



Use the samples from both chains to form an estimator

Proof

$$H \approx h(u^{(t_0)}) + \sum_{t=t_0+1}^{\tau-1} \left(h(u^{(t)}) - h(\bar{u}^{(t-1)}) \right)$$



 $\tau = \inf\{t \ge 1 : u^{(t)} = \bar{u}^{(t-1)}\}$

$$\begin{split} H &= \mathbb{E}_{\pi(u)}[h(u)] = \mathbb{E}[h(u^{(t_0)})] + \sum_{t=t_0+1}^{\infty} \left(\mathbb{E}[h(u^{(t)})] - \mathbb{E}[h(u^{(t-1)})] \right) & \text{telescoping sum} \\ &= \mathbb{E}[h(u^{(t_0)})] + \sum_{t=t_0+1}^{\infty} \left(\mathbb{E}[h(u^{(t)})] - \mathbb{E}[h(\overline{u}^{(t-1)})] \right) & \text{same marginals} \\ &= \mathbb{E}\left[h(u^{(t_0)}) + \sum_{t=t_0+1}^{\infty} \left(h(u^{(t)}) - h(\overline{u}^{(t-1)}) \right) \right] & \text{swap expectation and limit} \\ &= \mathbb{E}\left[h(u^{(t_0)}) + \sum_{t=t_0+1}^{\tau-1} \left(h(u^{(t)}) - h(\overline{u}^{(t-1)}) \right) + \sum_{t=\tau}^{\infty} \left(h(u^{(t)}) - h(\overline{u}^{(t-1)}) \right) \right] & \text{by design} \end{split}$$





DeepMind

Unbiased Estimators on an Extended Space

Our Proposal

- Start with $\nabla_{\theta} \mathcal{L}(\theta) = \mathbb{E}_{p_{\theta}(z \mid x)} \left[\nabla_{\theta} \log p_{\theta}(x, z) \right]$
- Form the augmented model and augmented proposal

$$p_{\theta,\phi}(x, z_{1:K}, \ell) = \frac{1}{K} p_{\theta}(z_{\ell}) p_{\theta}(x \mid z_{\ell}) \prod_{\substack{k=1\\k \neq \ell}}^{K} q_{\phi}(z_k \mid x)$$
$$q_{\theta,\phi}(z_{1:K}, \ell) = \text{Categorical}\left(\ell \mid \widetilde{w}_{\theta,\phi}^{(1)}, \dots, \widetilde{w}_{\theta,\phi}^{(K)}\right) \prod_{k=1}^{K} q_{\phi}(z_k \mid x)$$



 $egin{aligned} u &= [z_{1:K}, \ell] \ ar u &= [ar z_{1:K}, ar \ell] \end{aligned}$

- Run a coupled MCMC kernel on the extended space, targeting the augmented posterior
 - How to form the kernel?



Algorithm 1: Particle independent Metropolis-Hastings (PIMH) kernel, $\mathcal{K}_{\text{PIMH}}(\cdot, \cdot | z_{1:K}, \ell)$ Input: Current state of the chain, $(z_{1:K}, \ell)$ Output: Next state of the chain1 Sample a candidate $(z_{1:K}^*, \ell^*) \sim q_{\theta,\phi}(\cdot, \cdot)$ 2 Sample $u \sim \mathcal{U}([0, 1])$ 3 if $u \leq \alpha(z_{1:K}^*, \ell^* | z_{1:K}, \ell)$ then4 | Return $(z_{1:K}^*, \ell^*)$ > the proposal is accepted5 else6 | Return $(z_{1:K}, \ell)$ > the proposal is rejected7 end





Algorithm 4: Coupled PIMH kernel, $\mathcal{K}_{C-PIMH}((\cdot, \cdot), (\cdot, \cdot) | (z_{1:K}, \ell), (\bar{z}_{1:K}, \ell))$ **Input:** Current state of both chains, $(z_{1:K}, \ell)$ and $(\bar{z}_{1:K}, \ell)$ **Output:** New state of both chains 1 Sample $(z_{1\cdot K}^{\star}, \ell^{\star}) \sim q_{\theta,\phi}(\cdot, \cdot)$ 2 Sample $u \sim \mathcal{U}([0, 1])$ **3** if $u \leq \alpha(z_{1\cdot K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u \leq \alpha(z_{1\cdot K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \bar{\ell})$ then 4 Return $((z_{1\cdot K}^{\star}, \ell^{\star}), (z_{1\cdot K}^{\star}, \ell^{\star}))$ \triangleright both chains accept the proposal 5 else if $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u > \alpha(z_{1:K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \ell)$ then 6 Return $((z_{1:K}^{\star}, \ell^{\star}), (\bar{z}_{1:K}, \ell))$ \triangleright the first chain accepts the proposal 7 else if $u > \alpha(z_{1:K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \bar{\ell})$ then Return $((z_{1:K}, \ell), (z_{1:K}^{\star}, \ell^{\star}))$ 8 \triangleright the second chain accepts the proposal 9 else Return $((z_{1:K}, \ell), (\bar{z}_{1:K}, \ell))$ \triangleright neither chain accepts the proposal 10 11 end



Algorithm 4: Coupled PIMH kernel, $\mathcal{K}_{C-PIMH}((\cdot, \cdot), (\cdot, \cdot) | (z_{1:K}, \ell), (\bar{z}_{1:K}, \ell))$ **Input:** Current state of both chains, $(z_{1:K}, \ell)$ and $(\overline{z}_{1:K}, \overline{\ell})$ **Output:** New state of both chains 1 Sample $(z_{1\cdot K}^{\star}, \ell^{\star}) \sim q_{\theta,\phi}(\cdot, \cdot)$ 2 Sample $u \sim \mathcal{U}([0, 1])$ **3 if** $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \bar{\ell})$ then Return $((z_{1\cdot K}^{\star}, \ell^{\star}), (z_{1\cdot K}^{\star}, \ell^{\star}))$ 4 \triangleright both chains accept the proposal 5 else if $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u > \alpha(z_{1:K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \bar{\ell})$ then 6 Return $((z_{1\cdot K}^{\star}, \ell^{\star}), (\bar{z}_{1:K}, \ell))$ \triangleright the first chain accepts the proposal 7 else if $u > \alpha(z_{1:K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \bar{\ell})$ then Return $((z_{1:K}, \ell), (z_{1:K}^{\star}, \ell^{\star}))$ ▷ the second chain accepts the proposal 8 9 else Return $((z_{1:K}, \ell), (\bar{z}_{1:K}, \bar{\ell}))$ \triangleright neither chain accepts the proposal 10 11 end



Algorithm 4: Coupled PIMH kernel, $\mathcal{K}_{C-PIMH}((\cdot, \cdot), (\cdot, \cdot) | (z_{1:K}, \ell), (\bar{z}_{1:K}, \ell))$ **Input:** Current state of both chains, $(z_{1:K}, \ell)$ and $(\overline{z}_{1:K}, \ell)$ **Output:** New state of both chains 1 Sample $(z_{1\cdot K}^{\star}, \ell^{\star}) \sim q_{\theta,\phi}(\cdot, \cdot)$ 2 Sample $u \sim \mathcal{U}([0, 1])$ 3 if $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \bar{\ell})$ then Return $((z_{1\cdot K}^{\star}, \ell^{\star}), (z_{1\cdot K}^{\star}, \ell^{\star}))$ 4 \triangleright both chains accept the proposal 5 else if $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} | z_{1:K}, \ell)$ and $u > \alpha(z_{1:K}^{\star}, \ell^{\star} | \bar{z}_{1:K}, \bar{\ell})$ then 6 Return $((z_{1:K}^{\star}, \ell^{\star}), (\bar{z}_{1:K}, \ell))$ \triangleright the first chain accepts the proposal 7 else if $u > \alpha(z_{1:K}^{\star}, \ell^{\star} \mid z_{1:K}, \ell)$ and $u \leq \alpha(z_{1:K}^{\star}, \ell^{\star} \mid \overline{z}_{1:K}, \overline{\ell})$ then Return $((z_{1:K}, \ell), (z_{1:K}^{\star}, \ell^{\star}))$ \triangleright the second chain accepts the proposal 8 9 else Return $((z_{1:K}, \ell), (\bar{z}_{1:K}, \ell))$ \triangleright neither chain accepts the proposal 10 11 end



• After collecting samples, obtain the unbiased gradient estimator as

$$\nabla_{\theta} \log p_{\theta}(x) \approx h(z_{1:K}^{(t_0)}) + \sum_{t=t_0+1}^{\tau-1} \left(h(z_{1:K}^{(t)}) - h(\bar{z}_{1:K}^{(t-1)}) \right)$$

The function *h* is

$$h(z_{1:K}) = \sum_{k=1}^{K} \tilde{w}_{ heta,\phi}^{(k)}
abla_{ heta} \log p_{ heta}(x, z_k)$$



Our Contributions

- Our MCMC algorithm is **based on ISIR** (rather than PIMH)
- We propose an extension of ISIR, called DISIR, that significantly reduces the estimator variance
- We derive **sufficient conditions** that guarantee a finite-variance unbiased estimator in finite time
- Our estimator is based on a **lagged coupling estimator**, which further reduces the variance



Importance Sampling in High-Dimensional Spaces

- IS typically fails in high dimensions, when one weight dominates the others
- We augment the dimensionality with *K*-1 particles
 - So an IS-based MCMC algorithm should perform poorly (and the MCMC chains would never meet)
 - However, **performance improves with dimensionality** (and meeting occurs earlier) as the model and proposals become closer to each other when *K* increases









DeepMind

Experiments and Results

PPCA: Analysis of Unbiasedness





VAE on Binarized MNIST

train log-likelihood



Private & Confidential

Analysis of the Meeting Time





VAE on Fashion-MNIST and CIFAR-10





Conclusions

- The combination of latent space augmentation and coupling estimators gives practical unbiased gradients
 - Unbiased gradient estimation **improves the model's predictive performance** for VAEs
 - The **computational time is higher**, but we can use this method to refine model fits
- \rightarrow

 \rightarrow

 \rightarrow

Future work on improving coupling estimators will also reduce the computational complexity

