Challenges and Opportunities in Scalable Alpha-divergence Variational Inference: Application to IWAEs

Kamélia Daudel



 ${\sf CIRM}\,-\,{\sf Fusion}\,\,{\sf workshop}$ 

Joint work with Joe Benton, Arnaud Doucet and Yuyang Shi

# Outline

#### 1 Introduction

- 2 The VR-IWAE bound
- 3 Theoretical study of the VR-IWAE bound
- **4** Numerical experiments



# Outline

#### 1 Introduction

- 2 The VR-IWAE bound
- 3 Theoretical study of the VR-IWAE bound
- **4** Numerical experiments
- **5** Conclusion

#### • Setting :

- **(1)** We consider a model with joint distribution  $p_{\theta}(x, z)$  parameterized by  $\theta$ , where x is an observation and z is a latent variable valued in  $\mathbb{R}^d$
- ② In that case, the marginal log-likelihood of x is given by

$$\ell(\theta; x) := \log p_{\theta}(x) = \log \left( \int p_{\theta}(x, z) dz \right)$$

• Goal : find  $\theta$  which best describes the observation x

$$\theta^{\star} = \operatorname{argmax}_{\theta} \, \ell(\theta; x)$$

(more generally  $\theta^{\star} = \operatorname{argmax}_{\theta} \sum_{i=1}^{T} \ell(\theta; x_i)$ )

- Setting :
  - **()** We consider a model with joint distribution  $p_{\theta}(x, z)$  parameterized by  $\theta$ , where x is an observation and z is a latent variable valued in  $\mathbb{R}^d$

**2** In that case, the marginal log-likelihood of x is given by

$$\ell(\theta; x) := \log p_{\theta}(x) = \log \left( \int p_{\theta}(x, z) \mathrm{d}z \right)$$

• Goal : find  $\theta$  which best describes the observation x

$$\theta^{\star} = \operatorname{argmax}_{\theta} \, \ell(\theta; x)$$

(more generally  $\theta^{\star} = \operatorname{argmax}_{\theta} \sum_{i=1}^{T} \ell(\theta; x_i)$ )

- Setting :
  - **1** We consider a model with joint distribution  $p_{\theta}(x, z)$  parameterized by  $\theta$ , where x is an observation and z is a latent variable valued in  $\mathbb{R}^d$
  - **2** In that case, the marginal log-likelihood of x is given by

$$\ell(\theta; x) := \log p_{\theta}(x) = \log \left( \int p_{\theta}(x, z) dz \right)$$

• Goal : find  $\theta$  which best describes the observation x

$$\theta^{\star} = \operatorname{argmax}_{\theta} \ell(\theta; x)$$

(more generally  $\theta^{\star} = \operatorname{argmax}_{\theta} \sum_{i=1}^{T} \ell(\theta; x_i)$ )

- Setting :
  - **1** We consider a model with joint distribution  $p_{\theta}(x, z)$  parameterized by  $\theta$ , where x is an observation and z is a latent variable valued in  $\mathbb{R}^d$
  - **2** In that case, the marginal log-likelihood of x is given by

$$\ell(\theta; x) := \log p_{\theta}(x) = \log \left( \int p_{\theta}(x, z) dz \right)$$

• Goal : find  $\theta$  which best describes the observation x

 $\theta^* = \operatorname{argmax}_{\theta} \ell(\theta; x)$ 

(more generally  $\theta^{\star} = \operatorname{argmax}_{\theta} \sum_{i=1}^{T} \ell(\theta; x_i)$ )

- Setting :
  - **1** We consider a model with joint distribution  $p_{\theta}(x, z)$  parameterized by  $\theta$ , where x is an observation and z is a latent variable valued in  $\mathbb{R}^d$
  - **2** In that case, the marginal log-likelihood of x is given by

$$\ell(\theta; x) := \log p_{\theta}(x) = \log \left( \int p_{\theta}(x, z) dz \right)$$

• Goal : find  $\theta$  which best describes the observation x

$$\theta^{\star} = \operatorname{argmax}_{\theta} \, \ell(\theta; x)$$

(more generally  $\theta^{\star} = \operatorname{argmax}_{\theta} \sum_{i=1}^{T} \ell(\theta; x_i)$ )

- Setting :
  - **()** We consider a model with joint distribution  $p_{\theta}(x, z)$  parameterized by  $\theta$ , where x is an observation and z is a latent variable valued in  $\mathbb{R}^d$
  - **2** In that case, the marginal log-likelihood of x is given by

$$\ell(\theta; x) := \log p_{\theta}(x) = \log \left( \int p_{\theta}(x, z) dz \right)$$

• Goal : find  $\theta$  which best describes the observation x

$$\theta^{\star} = \operatorname{argmax}_{\theta} \ell(\theta; x)$$

(more generally  $\theta^{\star} = \operatorname{argmax}_{\theta} \sum_{i=1}^{T} \ell(\theta; x_i)$ )

- Setting :
  - **()** We consider a model with joint distribution  $p_{\theta}(x, z)$  parameterized by  $\theta$ , where x is an observation and z is a latent variable valued in  $\mathbb{R}^d$
  - **2** In that case, the marginal log-likelihood of x is given by

$$\ell(\theta; x) := \log p_{\theta}(x) = \log \left( \int p_{\theta}(x, z) dz \right)$$

• Goal : find  $\theta$  which best describes the observation x

$$\theta^{\star} = \operatorname{argmax}_{\theta} \, \ell(\theta; x)$$

(more generally  $\theta^{\star} = \operatorname{argmax}_{\theta} \sum_{i=1}^{T} \ell(\theta; x_i)$ )

- Idea : construct variational bounds, i.e. surrogate objective functions to the marginal log-likelihood that are more amenable to optimization.
- Common examples :

 $\rightarrow$  Evidence Lower BOund (ELBO) : rely on a variational probability density  $q_{\phi}(z|x)$  parameterized by  $\phi$ 

ELBO
$$(\theta, \phi; x) = \int q_{\phi}(z|x) \log(w_{\theta, \phi}(z; x)) dz$$
 where  $w_{\theta, \phi}(z; x) = \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)}$ 

with  $\operatorname{ELBO}(\theta,\phi;x) \leq \ell(\theta;x)$ 

ightarrow Importance Weighted Auto-Encoder (IWAE) bound (Burda et al., 2016)

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) = \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right) dz_{1:N}, \quad N \in \mathbb{N}^*$$

with  $\ell_N^{(\mathrm{IWAE})}(\theta,\phi;x) \leq \ell(\theta;x)$  and the unbiased Monte Carlo estimate

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) \approx \log\left(\frac{1}{N}\sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right), \quad z_j \sim q_\phi(\cdot|x), \quad j = 1\dots N$$

Kamélia Daudel (University of Oxford)  $\cdot$  Challenges and Opportunities in Scalable  $\alpha$ -divergence Variational Inference 5 / 37

- Idea : construct variational bounds, i.e. surrogate objective functions to the marginal log-likelihood that are more amenable to optimization.
- Common examples :

 $\rightarrow$  Evidence Lower BOund (ELBO) : rely on a variational probability density  $q_{\phi}(z|x)$  parameterized by  $\phi$ 

ELBO
$$(\theta, \phi; x) = \int q_{\phi}(z|x) \log(w_{\theta, \phi}(z; x)) dz$$
 where  $w_{\theta, \phi}(z; x) = \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)}$ 

with  $\operatorname{ELBO}(\theta,\phi;x) \leq \ell(\theta;x)$ 

ightarrow Importance Weighted Auto-Encoder (IWAE) bound (Burda et al., 2016)

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) = \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right) \mathrm{d}z_{1:N}, \quad N \in \mathbb{N}^*$$

with  $\ell_N^{(\mathrm{IWAE})}(\theta,\phi;x) \leq \ell(\theta;x)$  and the unbiased Monte Carlo estimate

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) \approx \log\left(\frac{1}{N}\sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right), \quad z_j \sim q_\phi(\cdot|x), \quad j = 1\dots N$$

Kamélia Daudel (University of Oxford)  $\cdot$  Challenges and Opportunities in Scalable  $\alpha$ -divergence Variational Inference 5 / 37

- Idea : construct variational bounds, i.e. surrogate objective functions to the marginal log-likelihood that are more amenable to optimization.
- Common examples :

 $\to$  Evidence Lower BOund (ELBO) : rely on a variational probability density  $q_\phi(z|x)$  parameterized by  $\phi$ 

$$\text{ELBO}(\theta,\phi;x) = \int q_{\phi}(z|x) \log\left(w_{\theta,\phi}(z;x)\right) \mathrm{d}z \quad \text{where} \quad w_{\theta,\phi}(z;x) = \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)}$$

with  $\mathrm{ELBO}(\theta,\phi;x) \leq \ell(\theta;x)$ 

ightarrow Importance Weighted Auto-Encoder (IWAE) bound (Burda et al., 2016)

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) = \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right) \mathrm{d}z_{1:N}, \quad N \in \mathbb{N}^*$$

with  $\ell_N^{(\mathrm{IWAE})}(\theta,\phi;x) \leq \ell(\theta;x)$  and the unbiased Monte Carlo estimate

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) \approx \log\left(\frac{1}{N}\sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right), \quad z_j \sim q_\phi(\cdot|x), \quad j = 1\dots N$$

Kamélia Daudel (University of Oxford)  $\cdot$  Challenges and Opportunities in Scalable  $\alpha$ -divergence Variational Inference 5 / 37

- Idea : construct variational bounds, i.e. surrogate objective functions to the marginal log-likelihood that are more amenable to optimization.
- Common examples :

 $\to$  Evidence Lower BOund (ELBO) : rely on a variational probability density  $q_\phi(z|x)$  parameterized by  $\phi$ 

$$\text{ELBO}(\theta,\phi;x) = \int q_{\phi}(z|x) \log\left(w_{\theta,\phi}(z;x)\right) dz \quad \text{where} \quad w_{\theta,\phi}(z;x) = \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)}$$

with  $\mathrm{ELBO}(\theta,\phi;x) \leq \ell(\theta;x)$ 

 $\rightarrow$  Importance Weighted Auto-Encoder (IWAE) bound (Burda et al., 2016)

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) = \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right) \mathrm{d}z_{1:N}, \quad N \in \mathbb{N}^\star$$

with  $\ell_N^{(\mathrm{IWAE})}(\theta,\phi;x) \leq \ell(\theta;x)$  and the unbiased Monte Carlo estimate

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) \approx \log\left(\frac{1}{N}\sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right), \quad z_j \sim q_\phi(\cdot|x), \quad j = 1\dots N$$

- Idea : construct variational bounds, i.e. surrogate objective functions to the marginal log-likelihood that are more amenable to optimization.
- Common examples :

 $\to$  Evidence Lower BOund (ELBO) : rely on a variational probability density  $q_\phi(z|x)$  parameterized by  $\phi$ 

$$\text{ELBO}(\theta,\phi;x) = \int q_{\phi}(z|x) \log\left(w_{\theta,\phi}(z;x)\right) dz \quad \text{where} \quad w_{\theta,\phi}(z;x) = \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)}$$

with  $\mathrm{ELBO}(\theta,\phi;x) \leq \ell(\theta;x)$ 

 $\rightarrow$  Importance Weighted Auto-Encoder (IWAE) bound (Burda et al., 2016)

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) = \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right) \mathrm{d}z_{1:N}, \quad N \in \mathbb{N}^\star$$

with  $\ell_N^{(\mathrm{IWAE})}(\theta,\phi;x) \leq \ell(\theta;x)$  and the unbiased Monte Carlo estimate

$$\ell_N^{(\text{IWAE})}(\theta,\phi;x) \approx \log\left(\frac{1}{N}\sum_{j=1}^N w_{\theta,\phi}(z_j;x)\right), \quad z_j \sim q_\phi(\cdot|x), \quad j = 1\dots N$$

Kamélia Daudel (University of Oxford) Challenges and Opportunities in Scalable lpha-divergence Variational Inference 5 / 37

# Training procedure for the IWAE bound

$${ig {?}}$$
 Reparameterization trick  $z=f(arepsilon,\phi;x)\sim q_\phi(\cdot|x)$  where  $arepsilon\sim q$ 

**Reparameterized** gradient estimator (Burda et al., 2016)

$$\frac{\partial}{\partial \phi} \ell_N^{(\text{IWAE})}(\theta, \phi; x) = \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j; x)}{\sum_{k=1}^N w_{\theta, \phi}(z_k; x)} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi; x); x) \right) \mathrm{d}\varepsilon_{1:N}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

ightarrow Unbiased SGD steps w.r.t.  $( heta,\phi)$ 

# Training procedure for the IWAE bound

$$\fbox{}$$
 Reparameterization trick  $z=f(arepsilon,\phi;x)\sim q_{\phi}(\cdot|x)$  where  $arepsilon\sim q$ 

Reparameterized gradient estimator (Burda et al., 2016)

$$\frac{\partial}{\partial \phi} \ell_N^{(\text{IWAE})}(\theta, \phi; x) = \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j; x)}{\sum_{k=1}^N w_{\theta, \phi}(z_k; x)} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi; x); x) \right) \mathrm{d}\varepsilon_{1:N}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

 $\rightarrow$  Unbiased SGD steps w.r.t.  $(\theta, \phi)$ 

# Training procedure for the IWAE bound

$$\fbox{}$$
 Reparameterization trick  $z=f(arepsilon,\phi;x)\sim q_{\phi}(\cdot|x)$  where  $arepsilon\sim q$ 

Reparameterized gradient estimator (Burda et al., 2016)

$$\frac{\partial}{\partial \phi} \ell_N^{(\text{IWAE})}(\theta, \phi; x) = \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j; x)}{\sum_{k=1}^N w_{\theta, \phi}(z_k; x)} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi; x); x) \right) \mathrm{d}\varepsilon_{1:N}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

 $\rightarrow$  Unbiased SGD steps w.r.t.  $(\theta, \phi)$ 

The Variational Rényi (VR) bound (Li and Turner, 2016) : for all  $\alpha \in \mathbb{R} \setminus \{1\}$ ,

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

ightarrow Lower bound on  $\ell( heta;x)$  for lpha>0 (upper for lpha<0)

 $\rightarrow$  Flexible family of variational bounds indexed by  $\alpha$  which recovers the ELBO when  $\alpha \rightarrow 1$  (also has ties to the  $\alpha$ -divergence)

Training procedure using the reparameterized gradient estimator

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta, \phi; x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z; x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon, \phi; x); x) \ d\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z; x)^{1-\alpha} \ d\varepsilon}$$
$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j; x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k; x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j, \phi; x); x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

The Variational Rényi (VR) bound (Li and Turner, 2016) : for all  $\alpha \in \mathbb{R} \setminus \{1\}$ ,

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

 $\rightarrow$  Lower bound on  $\ell(\theta; x)$  for  $\alpha > 0$  (upper for  $\alpha < 0$ )

 $\rightarrow$  Flexible family of variational bounds indexed by  $\alpha$  which recovers the ELBO when  $\alpha \rightarrow 1$  (also has ties to the  $\alpha$ -divergence)

Training procedure using the reparameterized gradient estimator

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta, \phi; x) = \frac{\int q(\varepsilon) \ w_{\theta, \phi}(z; x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon, \phi; x); x) \ d\varepsilon}{\int q(\varepsilon) \ w_{\theta, \phi}(z; x)^{1-\alpha} \ d\varepsilon}$$
$$\approx \sum_{j=1}^{N} \frac{w_{\theta, \phi}(z_j; x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta, \phi}(z_k; x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi; x); x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

The Variational Rényi (VR) bound (Li and Turner, 2016) : for all  $\alpha \in \mathbb{R} \setminus \{1\}$ ,

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

 $\rightarrow$  Lower bound on  $\ell(\theta; x)$  for  $\alpha > 0$  (upper for  $\alpha < 0$ )

 $\rightarrow$  Flexible family of variational bounds indexed by  $\alpha$  which recovers the ELBO when  $\alpha \rightarrow 1$  (also has ties to the  $\alpha$ -divergence)

Training procedure using the reparameterized gradient estimator

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon,\phi;x);x) \ d\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ d\varepsilon}$$
$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

The Variational Rényi (VR) bound (Li and Turner, 2016) : for all  $\alpha \in \mathbb{R} \setminus \{1\}$ ,

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

 $\rightarrow$  Lower bound on  $\ell(\theta; x)$  for  $\alpha > 0$  (upper for  $\alpha < 0$ )

 $\rightarrow$  Flexible family of variational bounds indexed by  $\alpha$  which recovers the ELBO when  $\alpha \rightarrow 1$  (also has ties to the  $\alpha$ -divergence)

Training procedure using the reparameterized gradient estimator

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta, \phi; x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z; x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon, \phi; x); x) \ d\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z; x)^{1-\alpha} \ d\varepsilon}$$
$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j; x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k; x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j, \phi; x); x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

The Variational Rényi (VR) bound (Li and Turner, 2016) : for all  $\alpha \in \mathbb{R} \setminus \{1\}$ ,

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

 $\rightarrow$  Lower bound on  $\ell(\theta; x)$  for  $\alpha > 0$  (upper for  $\alpha < 0$ )

 $\rightarrow$  Flexible family of variational bounds indexed by  $\alpha$  which recovers the ELBO when  $\alpha \rightarrow 1$  (also has ties to the  $\alpha$ -divergence)

Training procedure using the reparameterized gradient estimator

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon,\phi;x);x) \ d\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ d\varepsilon}$$
$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

with positive empirical results

 $\Im$  Recovers SGD with the IWAE bound for  $\alpha = 0$  (resp. ELBO for  $\alpha = 1$ )

Kamélia Daudel (University of Oxford)  $\cdot$  Challenges and Opportunities in Scalable  $\alpha$ -divergence Variational Inference 7 / 37

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon,\phi;x);x) \ \mathrm{d}\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}\varepsilon}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

- $\rightarrow$  However
  - The VR bound can only be estimated using biased MC estimators
  - **②** The VR bound does not recover the IWAE bound when  $\alpha = 0$
  - **(a)** No theoretical justification as SGD with the VR bound resorts to biased estimators on top of the reparameterization trick (unless  $\alpha \in \{0, 1\}$ )

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon,\phi;x);x) \ \mathrm{d}\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}\varepsilon}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

#### $\rightarrow$ However

• The VR bound can only be estimated using biased MC estimators

- ② The VR bound does not recover the IWAE bound when  $\alpha = 0$
- **(a)** No theoretical justification as SGD with the VR bound resorts to biased estimators on top of the reparameterization trick (unless  $\alpha \in \{0, 1\}$ )

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon,\phi;x);x) \ \mathrm{d}\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}\varepsilon}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

#### $\rightarrow$ However

- The VR bound can only be estimated using biased MC estimators
- ② The VR bound does not recover the IWAE bound when  $\alpha = 0$
- **(a)** No theoretical justification as SGD with the VR bound resorts to biased estimators on top of the reparameterization trick (unless  $\alpha \in \{0, 1\}$ )

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon,\phi;x);x) \ \mathrm{d}\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}\varepsilon}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

- $\rightarrow$  However
  - The VR bound can only be estimated using biased MC estimators
  - 2 The VR bound does not recover the IWAE bound when  $\alpha=0$
  - **③** No theoretical justification as SGD with the VR bound resorts to biased estimators on top of the reparameterization trick (unless  $\alpha \in \{0, 1\}$ )

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(\cdot|x), \quad j = 1 \dots N$$

$$\frac{\partial}{\partial \phi} \mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon,\phi;x);x) \ \mathrm{d}\varepsilon}{\int q(\varepsilon) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}\varepsilon}$$

$$\approx \sum_{j=1}^{N} \frac{w_{\theta,\phi}(z_j;x)^{1-\alpha}}{\sum_{k=1}^{N} w_{\theta,\phi}(z_k;x)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi;x);x), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

- $\rightarrow$  However
  - The VR bound can only be estimated using biased MC estimators
  - 2 The VR bound does not recover the IWAE bound when  $\alpha=0$
  - **③** No theoretical justification as SGD with the VR bound resorts to biased estimators on top of the reparameterization trick (unless  $\alpha \in \{0, 1\}$ )

• Li and Turner (Theorem 2, 2016) further looked into the biased approximation of the VR bound

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(z|x), \quad j = 1 \dots N$$

• They investigated the expectation of the biased MC approximation, i.e.

$$\frac{1}{1-\alpha} \int \int \prod_{i=1}^{N} q_{\phi}(z_i|x) \log \left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

e.g. they showed that it is non-decreasing with N when  $\alpha \leq 1$ .

• Question Could this expectation be seen as a variational bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• Li and Turner (Theorem 2, 2016) further looked into the biased approximation of the VR bound

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(z|x), \quad j = 1 \dots N$$

• They investigated the expectation of the biased MC approximation, i.e.

$$\frac{1}{1-\alpha} \int \int \prod_{i=1}^{N} q_{\phi}(z_i|x) \log \left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

e.g. they showed that it is non-decreasing with N when  $\alpha \leq 1$ .

• <u>Question</u> Could this expectation be seen as a variational bound? Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• Li and Turner (Theorem 2, 2016) further looked into the biased approximation of the VR bound

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(z|x), \quad j = 1 \dots N$$

• They investigated the expectation of the biased MC approximation, i.e.

$$\frac{1}{1-\alpha} \int \int \prod_{i=1}^{N} q_{\phi}(z_i|x) \log \left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

e.g. they showed that it is non-decreasing with N when  $\alpha \leq 1.$ 

• Question Could this expectation be seen as a variational bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• Li and Turner (Theorem 2, 2016) further looked into the biased approximation of the VR bound

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) = \frac{1}{1-\alpha} \log\left(\int q_{\phi}(z|x) \ w_{\theta,\phi}(z;x)^{1-\alpha} \ \mathrm{d}z\right)$$
$$\approx \frac{1}{1-\alpha} \log\left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right), \quad z_j \sim q_{\phi}(z|x), \quad j = 1 \dots N$$

• They investigated the expectation of the biased MC approximation, i.e.

$$\frac{1}{1-\alpha} \int \int \prod_{i=1}^{N} q_{\phi}(z_i|x) \log \left(\frac{1}{N} \sum_{j=1}^{N} w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

e.g. they showed that it is non-decreasing with N when  $\alpha \leq 1.$ 

• Question Could this expectation be seen as a variational bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

# Outline

### 1 Introduction

### 2 The VR-IWAE bound

#### 3 Theoretical study of the VR-IWAE bound

#### **4** Numerical experiments

#### **5** Conclusion

# The VR-IWAE bound

For all  $\alpha \in [0,1)$  and all  $N \in \mathbb{N}^{\star}$ 

$$\ell_N^{(\alpha)}(\theta,\phi;x) := \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

The VR-IWAE bound is a **lower bound** on the marginal log-likelihood that

- 1 Can be estimated using unbiased MC estimators
- **2** Recovers the IWAE objective function when  $\alpha = 0$
- Scale Leads to the same SGD procedure as the VR bound in the reparameterized case, but this time using unbiased estimators

$$\begin{split} \frac{\partial}{\partial \phi} \ell_N^{(\alpha)}(\theta, \phi; x) \\ &= \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta, \phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi)) \right) d\varepsilon_{1:N}. \\ &\approx \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta, \phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi)), \quad \varepsilon_j \sim q, \quad j = 1 \dots N \end{split}$$

# The VR-IWAE bound

For all  $\alpha \in [0,1)$  and all  $N \in \mathbb{N}^{\star}$ 

$$\ell_N^{(\alpha)}(\theta,\phi;x) := \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

The VR-IWAE bound is a lower bound on the marginal log-likelihood that

- 1 Can be estimated using unbiased MC estimators
- **2** Recovers the IWAE objective function when  $\alpha = 0$
- Scale Leads to the same SGD procedure as the VR bound in the reparameterized case, but this time using unbiased estimators

$$\begin{split} \frac{\partial}{\partial \phi} \ell_N^{(\alpha)}(\theta,\phi;x) \\ &= \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi)) \right) d\varepsilon_{1:N}. \\ &\approx \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi)), \quad \varepsilon_j \sim q, \quad j = 1 \dots N \end{split}$$

# The VR-IWAE bound

For all  $\alpha \in [0,1)$  and all  $N \in \mathbb{N}^{\star}$ 

$$\ell_N^{(\alpha)}(\theta,\phi;x) := \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

The VR-IWAE bound is a lower bound on the marginal log-likelihood that

Can be estimated using unbiased MC estimators

- **2** Recovers the IWAE objective function when  $\alpha = 0$
- Searching Leads to the same SGD procedure as the VR bound in the reparameterized case, but this time using unbiased estimators

$$\begin{split} \frac{\partial}{\partial \phi} \ell_N^{(\alpha)}(\theta,\phi;x) \\ &= \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi)) \right) d\varepsilon_{1:N}. \\ &\approx \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi)), \quad \varepsilon_j \sim q, \quad j = 1 \dots N \end{split}$$
For all  $\alpha \in [0,1)$  and all  $N \in \mathbb{N}^{\star}$ 

$$\ell_N^{(\alpha)}(\theta,\phi;x) := \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

- Can be estimated using unbiased MC estimators
- **2** Recovers the IWAE objective function when  $\alpha=0$
- Scale Leads to the same SGD procedure as the VR bound in the reparameterized case, but this time using unbiased estimators

$$\begin{split} \frac{\partial}{\partial \phi} \ell_N^{(\alpha)}(\theta, \phi; x) \\ &= \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta, \phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi)) \right) d\varepsilon_{1:N}. \\ &\approx \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta, \phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi)), \quad \varepsilon_j \sim q, \quad j = 1 \dots N \end{split}$$

For all  $\alpha \in [0,1)$  and all  $N \in \mathbb{N}^{\star}$ 

$$\ell_N^{(\alpha)}(\theta,\phi;x) := \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

- ① Can be estimated using unbiased MC estimators
- **2** Recovers the IWAE objective function when  $\alpha = 0$
- Leads to the same SGD procedure as the VR bound in the reparameterized case, but this time using unbiased estimators

$$\frac{\partial}{\partial \phi} \ell_N^{(\alpha)}(\theta, \phi; x) = \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j, \phi)) \right) d\varepsilon_{1:N}.$$
$$\approx \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j, \phi)), \quad \varepsilon_j \sim q, \quad j = 1 \dots N$$

For all  $\alpha \in [0,1)$  and all  $N \in \mathbb{N}^{\star}$ 

$$\ell_N^{(\alpha)}(\theta,\phi;x) := \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

- Can be estimated using unbiased MC estimators
- **2** Recovers the IWAE objective function when  $\alpha = 0$
- Leads to the same SGD procedure as the VR bound in the reparameterized case, but this time using unbiased estimators

$$\begin{split} &\frac{\partial}{\partial \phi} \ell_N^{(\alpha)}(\theta,\phi;x) \\ &= \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi)) \right) \mathrm{d}\varepsilon_{1:N}. \\ &\approx \sum_{j=1}^N \frac{w_{\theta,\phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta,\phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta,\phi}(f(\varepsilon_j,\phi)), \quad \varepsilon_j \sim q, \quad j = 1 \dots N \end{split}$$

For all  $\alpha \in [0,1)$  and all  $N \in \mathbb{N}^{\star}$ 

$$\ell_N^{(\alpha)}(\theta,\phi;x) := \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N w_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

- Can be estimated using unbiased MC estimators
- **2** Recovers the IWAE objective function when  $\alpha = 0$
- Leads to the same SGD procedure as the VR bound in the reparameterized case, but this time using unbiased estimators

$$\begin{split} &\frac{\partial}{\partial \phi} \ell_N^{(\alpha)}(\theta, \phi; x) \\ &= \int \int \prod_{i=1}^N q(\varepsilon_i) \left( \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta, \phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi)) \right) \mathrm{d}\varepsilon_{1:N}. \\ &\approx \sum_{j=1}^N \frac{w_{\theta, \phi}(z_j)^{1-\alpha}}{\sum_{k=1}^N w_{\theta, \phi}(z_k)^{1-\alpha}} \frac{\partial}{\partial \phi} \log w_{\theta, \phi}(f(\varepsilon_j, \phi)), \quad \varepsilon_j \sim q, \quad j = 1 \dots N \end{split}$$

 $\rightarrow$  The VR-IWAE bound is the **theoretically-sound** extension of the IWAE bound originating from the  $\alpha\text{-divergence VI}$  methodology

 $\rightarrow$  It provides **theoretical guarantees** behind various VR-bound gradient-based schemes previously proposed in the  $\alpha$ -divergence VI community

ightarrow Other notable advantages of the VR-IWAE bound :

$$\begin{split} & \mathrm{SNR}_{\theta_{\ell}} = \Theta(\sqrt{N}) \\ & \mathrm{SNR}_{\phi_{\ell'}} = \begin{cases} \Theta(\sqrt{1/N}) & \text{if } \alpha = 0 \text{ (Rainforth et al., 2018),} \\ & \Theta(\sqrt{N}) & \text{if } \alpha \in (0, 1). \end{cases} \end{split}$$

- The doubly-reparameterized gradient estimators of the IWAE generalize to the VR-IWAE bound
- ightarrow Motivates the use of  $lpha \in [0,1)$  in practice

 $\rightarrow$  The VR-IWAE bound is the **theoretically-sound** extension of the IWAE bound originating from the  $\alpha$ -divergence VI methodology

 $\rightarrow$  It provides **theoretical guarantees** behind various VR-bound gradient-based schemes previously proposed in the  $\alpha$ -divergence VI community

ightarrow Other notable advantages of the VR-IWAE bound :

$$\begin{split} & \mathrm{SNR}_{\theta_{\ell}} = \Theta(\sqrt{N}) \\ & \mathrm{SNR}_{\phi_{\ell'}} = \begin{cases} \Theta(\sqrt{1/N}) & \text{if } \alpha = 0 \text{ (Rainforth et al., 2018),} \\ & \Theta(\sqrt{N}) & \text{if } \alpha \in (0, 1). \end{cases} \end{split}$$

- The doubly-reparameterized gradient estimators of the IWAE generalize to the VR-IWAE bound
- ightarrow Motivates the use of  $lpha \in [0,1)$  in practice

 $\rightarrow$  The VR-IWAE bound is the **theoretically-sound** extension of the IWAE bound originating from the  $\alpha$ -divergence VI methodology

 $\rightarrow$  It provides **theoretical guarantees** behind various VR-bound gradient-based schemes previously proposed in the  $\alpha$ -divergence VI community

 $\rightarrow$  Other notable advantages of the VR-IWAE bound :

$$\begin{split} & \mathrm{SNR}_{\theta_\ell} = \Theta(\sqrt{N}) \\ & \mathrm{SNR}_{\phi_{\ell'}} = \begin{cases} \Theta(\sqrt{1/N}) & \text{if } \alpha = 0 \text{ (Rainforth et al., 2018),} \\ & \Theta(\sqrt{N}) & \text{if } \alpha \in (0, 1). \end{cases} \end{split}$$

- The doubly-reparameterized gradient estimators of the IWAE generalize to the VR-IWAE bound
- ightarrow Motivates the use of  $lpha \in [0,1)$  in practice

 $\rightarrow$  The VR-IWAE bound is the **theoretically-sound** extension of the IWAE bound originating from the  $\alpha$ -divergence VI methodology

 $\rightarrow$  It provides **theoretical guarantees** behind various VR-bound gradient-based schemes previously proposed in the  $\alpha$ -divergence VI community

 $\rightarrow$  Other notable advantages of the VR-IWAE bound :

$$\begin{aligned} & \mathrm{SNR}_{\theta_{\ell}} = \Theta(\sqrt{N}) \\ & \mathrm{SNR}_{\phi_{\ell'}} = \begin{cases} \Theta(\sqrt{1/N}) & \text{if } \alpha = 0 \text{ (Rainforth et al., 2018),} \\ & \Theta(\sqrt{N}) & \text{if } \alpha \in (0, 1). \end{cases} \end{aligned}$$

- The doubly-reparameterized gradient estimators of the IWAE generalize to the VR-IWAE bound
- ightarrow Motivates the use of  $lpha \in [0,1)$  in practice

 $\rightarrow$  The VR-IWAE bound is the **theoretically-sound** extension of the IWAE bound originating from the  $\alpha\text{-divergence VI}$  methodology

 $\rightarrow$  It provides **theoretical guarantees** behind various VR-bound gradient-based schemes previously proposed in the  $\alpha$ -divergence VI community

 $\rightarrow$  Other notable advantages of the VR-IWAE bound :

$$\begin{aligned} & \mathrm{SNR}_{\theta_{\ell}} = \Theta(\sqrt{N}) \\ & \mathrm{SNR}_{\phi_{\ell'}} = \begin{cases} \Theta(\sqrt{1/N}) & \text{if } \alpha = 0 \text{ (Rainforth et al., 2018),} \\ & \Theta(\sqrt{N}) & \text{if } \alpha \in (0, 1). \end{cases} \end{aligned}$$

- The doubly-reparameterized gradient estimators of the IWAE generalize to the VR-IWAE bound
- ightarrow Motivates the use of  $lpha \in [0,1)$  in practice

 $\rightarrow$  The VR-IWAE bound is the **theoretically-sound** extension of the IWAE bound originating from the  $\alpha\text{-divergence VI}$  methodology

 $\rightarrow$  It provides **theoretical guarantees** behind various VR-bound gradient-based schemes previously proposed in the  $\alpha$ -divergence VI community

 $\rightarrow$  Other notable advantages of the VR-IWAE bound :

$$\begin{aligned} & \mathrm{SNR}_{\theta_{\ell}} = \Theta(\sqrt{N}) \\ & \mathrm{SNR}_{\phi_{\ell'}} = \begin{cases} \Theta(\sqrt{1/N}) & \text{if } \alpha = 0 \text{ (Rainforth et al., 2018),} \\ & \Theta(\sqrt{N}) & \text{if } \alpha \in (0, 1). \end{cases} \end{aligned}$$

- The doubly-reparameterized gradient estimators of the IWAE generalize to the VR-IWAE bound
- $\rightarrow$  Motivates the use of  $\alpha \in [0,1)$  in practice

## Outline

### 1 Introduction

### 2 The VR-IWAE bound

### 3 Theoretical study of the VR-IWAE bound

- **4** Numerical experiments
- **5** Conclusion

## Outline

### 1 Introduction

#### 2 The VR-IWAE bound

#### 3 Theoretical study of the VR-IWAE bound Overview First study

Second study



#### **5** Conclusion

 $\rightarrow$  Quantity of interest : variational gap

$$\begin{split} &\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1) \\ &= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N} \end{split}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z\in\mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)}$$

NB : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience  $\rightarrow$  Two **complentary** studies

 $\blacksquare$  When  $N \rightarrow \infty$  and the dimension of the latent space d is fixed

**2** When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\rightarrow$  Quantity of interest : variational gap

$$\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1)$$
$$= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z \in \mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)},$$

NB : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience  $\rightarrow$  Two **complentary** studies

 $\blacksquare$  When  $N \rightarrow \infty$  and the dimension of the latent space d is fixed

**2** When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\rightarrow$  Quantity of interest : variational gap

$$\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1)$$
$$= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z \in \mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)},$$

NB : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience

→ Two complentary studies

 $\blacksquare$  When  $N \rightarrow \infty$  and the dimension of the latent space d is fixed

2 When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\rightarrow$  Quantity of interest : variational gap

$$\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1)$$
$$= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z \in \mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)},$$

NB : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience

#### $\rightarrow$ Two complentary studies

() When  $N \to \infty$  and the dimension of the latent space d is fixed

**2** When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\rightarrow$  Quantity of interest : variational gap

$$\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1)$$
$$= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z \in \mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)},$$

NB : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience

#### $\rightarrow$ Two complentary studies

() When  $N \to \infty$  and the dimension of the latent space d is fixed

2 When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\rightarrow$  Quantity of interest : variational gap

$$\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1)$$
$$= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z \in \mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)},$$

NB : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience

#### $\rightarrow$ Two complentary studies

() When  $N \to \infty$  and the dimension of the latent space d is fixed

2 When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\rightarrow$  Quantity of interest : variational gap

$$\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1)$$
$$= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z \in \mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)},$$

NB : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience

#### $\rightarrow$ Two complentary studies

• When  $N \to \infty$  and the dimension of the latent space d is fixed • This analysis will be tailored for low to medium dimensions settings

**2** When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\rightarrow$  Quantity of interest : variational gap

$$\Delta_N^{(\alpha)}(\theta,\phi;x) := \ell_N^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x), \quad \alpha \in [0,1)$$
$$= \frac{1}{1-\alpha} \int \int \prod_{i=1}^N q_\phi(z_i|x) \log\left(\frac{1}{N} \sum_{j=1}^N \overline{w}_{\theta,\phi}(z_j;x)^{1-\alpha}\right) \mathrm{d}z_{1:N}$$

where  $\overline{w}_{\theta,\phi}(z_1;x),\ldots,\overline{w}_{\theta,\phi}(z_N;x)$  are the relative weights : for all  $z\in\mathbb{R}^d$ ,

$$\overline{w}_{\theta,\phi}(z;x) := \frac{w_{\theta,\phi}(z;x)}{\mathbb{E}_{Z \sim q_{\phi}}\left(w_{\theta,\phi}(Z;x)\right)} = \frac{w_{\theta,\phi}(z;x)}{p_{\theta}(x)} = \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)},$$

 $\mathsf{NB}$  : we will drop the dependency in x in  $\overline{w}_{\theta,\phi}(z;x)$  for convenience

#### $\rightarrow$ Two complentary studies

• When  $N \to \infty$  and the dimension of the latent space d is fixed • This analysis will be tailored for low to medium dimensions settings

**2** When 
$$N, d \to \infty$$
 with (i)  $\frac{\log N}{d} \to 0$  or (ii)  $\frac{\log N}{d^{1/3}} \to 0$ 

 $\ref{eq: the state of the state of the state of the state of the setting the setting the setting the state of the state o$ 

## Outline

### 1 Introduction

#### 2 The VR-IWAE bound

#### Theoretical study of the VR-IWAE bound Overview First study Second study

**4** Numerical experiments

#### **5** Conclusion

 $\rightarrow$  Maddison et al. (2017) followed by Domke and Sheldon (2018) looked into the variational gap for the IWAE bound ( $\alpha=0$ )

Informally, Domke and Sheldon (2018, Theorem 3) states that

$$\Delta_N^{(0)}(\theta,\phi;x) = -\frac{\gamma_0^2}{2N} + o\left(\frac{1}{N}\right)$$

where  $\gamma_0$  is the variance of the relative weights, i.e.

$$\gamma_0^2 := \mathbb{V}_{Z \sim q_\phi}(\overline{w}_{\theta,\phi}(Z))$$

- N is very beneficial to reduce  $\Delta_N^{(0)}( heta,\phi;x)$  (goes to 0 at a fast 1/N rate)
- Question What about  $\Delta_N^{(\alpha)}(\theta,\phi;x)$ ,  $\alpha \in [0,1)$ ?

 $\rightarrow$  Maddison et al. (2017) followed by Domke and Sheldon (2018) looked into the variational gap for the IWAE bound ( $\alpha=0$ )

Informally, Domke and Sheldon (2018, Theorem 3) states that

$$\Delta_N^{(0)}(\theta,\phi;x) = -\frac{\gamma_0^2}{2N} + o\left(\frac{1}{N}\right)$$

where  $\gamma_0$  is the variance of the relative weights, i.e.

$$\gamma_0^2 := \mathbb{V}_{Z \sim q_\phi}(\overline{w}_{\theta,\phi}(Z))$$

- N is very beneficial to reduce  $\Delta_N^{(0)}( heta,\phi;x)$  (goes to 0 at a fast 1/N rate)
- Question What about  $\Delta_N^{(\alpha)}(\theta,\phi;x)$ ,  $\alpha \in [0,1)$ ?

 $\rightarrow$  Maddison et al. (2017) followed by Domke and Sheldon (2018) looked into the variational gap for the IWAE bound ( $\alpha=0$ )

Informally, Domke and Sheldon (2018, Theorem 3) states that

$$\Delta_N^{(0)}(\theta,\phi;x) = -\frac{\gamma_0^2}{2N} + o\left(\frac{1}{N}\right)$$

where  $\gamma_0$  is the variance of the relative weights, i.e.

$$\gamma_0^2 := \mathbb{V}_{Z \sim q_\phi}(\overline{w}_{\theta,\phi}(Z))$$

- N is very beneficial to reduce  $\Delta_N^{(0)}( heta,\phi;x)$  (goes to 0 at a fast 1/N rate)
- Question What about  $\Delta_N^{(\alpha)}(\theta,\phi;x)$ ,  $\alpha \in [0,1)$ ?

## N goes to infinity and d is fixed

 $\rightarrow$  Maddison et al. (2017) followed by Domke and Sheldon (2018) looked into the variational gap for the IWAE bound ( $\alpha=0$ )

Informally, Domke and Sheldon (2018, Theorem 3) states that

$$\Delta_N^{(0)}(\theta,\phi;x) = -\frac{\gamma_0^2}{2N} + o\left(\frac{1}{N}\right)$$

where  $\gamma_0$  is the variance of the relative weights, i.e.

$$\gamma_0^2 := \mathbb{V}_{Z \sim q_\phi}(\overline{w}_{\theta,\phi}(Z))$$

- N is very beneficial to reduce  $\Delta_N^{(0)}(\theta,\phi;x)$  (goes to 0 at a fast 1/N rate)
- Question What about  $\Delta_N^{(\alpha)}(\theta,\phi;x)$ ,  $\alpha \in [0,1)$ ?

 $\rightarrow$  Maddison et al. (2017) followed by Domke and Sheldon (2018) looked into the variational gap for the IWAE bound ( $\alpha=0$ )

Informally, Domke and Sheldon (2018, Theorem 3) states that

$$\Delta_N^{(0)}(\theta,\phi;x) = -\frac{\gamma_0^2}{2N} + o\left(\frac{1}{N}\right)$$

where  $\gamma_0$  is the variance of the relative weights, i.e.

$$\gamma_0^2 := \mathbb{V}_{Z \sim q_\phi}(\overline{w}_{\theta,\phi}(Z))$$

- N is very beneficial to reduce  $\Delta_N^{(0)}(\theta,\phi;x)$  (goes to 0 at a fast 1/N rate)
- Question What about  $\Delta_N^{(\alpha)}(\theta,\phi;x)$ ,  $\alpha \in [0,1)$ ?

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right).$$

 $\rightarrow$  Two main terms :

- ${\rm I}{\rm I}$  A term going to zero at a fast 1/N rate that depends on  $\gamma^2_\alpha$
- **2** An error term  $\mathcal{L}^{(\alpha)}(\theta, \phi; x) \ell(\theta; x)$  [decreases away from 0 as  $\alpha$  increases]

- ightarrow "some conditions"
  - generalize the conditions from Domke and Sheldon (2018)
  - do not get more restrictive as lpha increases, motivates  $lpha \in (0,1)$
  - one of them controls  $\gamma^2_{lpha}$

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right).$$

 $\rightarrow$  Two main terms :

 ${\bf 0}$  A term going to zero at a fast 1/N rate that depends on  $\gamma_{\alpha}^2$ 

**2** An error term  $\mathcal{L}^{(\alpha)}(\theta, \phi; x) - \ell(\theta; x)$  [decreases away from 0 as  $\alpha$  increases]

- ightarrow "some conditions"
  - generalize the conditions from Domke and Sheldon (2018)
  - do not get more restrictive as  $\alpha$  increases, motivates  $\alpha \in (0,1)$
  - one of them controls  $\gamma^2_{lpha}$

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right).$$

 $\rightarrow$  Two main terms :

 $\textbf{0} \ \ {\rm A \ term \ going \ to \ zero \ at \ a \ fast \ } 1/N \ {\rm rate \ that \ depends \ on \ } \gamma_\alpha^2$ 

**2** An error term  $\mathcal{L}^{(\alpha)}(\theta, \phi; x) - \ell(\theta; x)$  [decreases away from 0 as  $\alpha$  increases]

- ightarrow "some conditions"
  - generalize the conditions from Domke and Sheldon (2018)
  - do not get more restrictive as  $\alpha$  increases, motivates  $\alpha \in (0,1)$
  - one of them controls  $\gamma^2_{lpha}$

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right).$$

 $\rightarrow$  Two main terms :

 $\textbf{0} \ \ {\rm A \ term \ going \ to \ zero \ at \ a \ fast \ } 1/N \ {\rm rate \ that \ depends \ on \ } \gamma_\alpha^2$ 

**2** An error term  $\mathcal{L}^{(\alpha)}(\theta, \phi; x) - \ell(\theta; x)$  [decreases away from 0 as  $\alpha$  increases]

- → "some conditions"
  - generalize the conditions from Domke and Sheldon (2018)
  - do not get more restrictive as lpha increases, motivates  $lpha \in (0,1)$
  - one of them controls  $\gamma^2_{lpha}$

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right).$$

 $\rightarrow$  Two main terms :

 $\textbf{0} \ \ {\rm A \ term \ going \ to \ zero \ at \ a \ fast \ } 1/N \ {\rm rate \ that \ depends \ on \ } \gamma_\alpha^2$ 

**2** An error term  $\mathcal{L}^{(\alpha)}(\theta, \phi; x) - \ell(\theta; x)$  [decreases away from 0 as  $\alpha$  increases]

- → "some conditions"
  - generalize the conditions from Domke and Sheldon (2018)
  - do not get more restrictive as  $\alpha$  increases, motivates  $\alpha \in (0,1)$
  - one of them controls  $\gamma^2_{lpha}$

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right).$$

 $\rightarrow$  Two main terms :

() A term going to zero at a fast 1/N rate that depends on  $\gamma_{lpha}^2$ 

**2** An error term  $\mathcal{L}^{(\alpha)}(\theta, \phi; x) - \ell(\theta; x)$  [decreases away from 0 as  $\alpha$  increases]

- → "some conditions"
  - generalize the conditions from Domke and Sheldon (2018)
  - do not get more restrictive as  $\alpha$  increases, motivates  $\alpha \in (0,1)$
  - one of them controls  $\gamma^2_{lpha}$

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right).$$

 $\rightarrow$  Two main terms :

() A term going to zero at a fast 1/N rate that depends on  $\gamma_{lpha}^2$ 

**2** An error term  $\mathcal{L}^{(\alpha)}(\theta, \phi; x) - \ell(\theta; x)$  [decreases away from 0 as  $\alpha$  increases]

- → "some conditions"
  - generalize the conditions from Domke and Sheldon (2018)
  - do not get more restrictive as  $\alpha$  increases, motivates  $\alpha \in (0,1)$
  - one of them controls  $\gamma_{\alpha}^2$

## Main result (cont'd)

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:  $\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_{\alpha}^2}{2N} + o\left(\frac{1}{N}\right).$ 

 $\rightarrow$  To the best of our knowledge, first result shedding light on how  $\alpha$  may play a role in the the success of  $\alpha$ -divergence VI.

 $\rightarrow$  Question Can we find some limitations to this approach?

## Main result (cont'd)

#### Theorem

Let  $\alpha \in [0,1)$ , denote  $\overline{w}_{\theta,\phi}^{(\alpha)}(z) = w_{\theta,\phi}(z)^{1-\alpha} / \mathbb{E}_{Z \sim q_{\phi}}(w_{\theta,\phi}(Z)^{1-\alpha})$  for all  $z \in \mathbb{R}^d$ and  $\gamma_{\alpha}^2 = (1-\alpha)^{-1} \mathbb{V}_{Z \sim q_{\phi}}(\overline{w}_{\theta,\phi}^{(\alpha)}(Z))$ . Then, under "some conditions", we have:  $\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_{\alpha}^2}{2N} + o\left(\frac{1}{N}\right).$ 

 $\rightarrow$  To the best of our knowledge, first result shedding light on how  $\alpha$  may play a role in the the success of  $\alpha$ -divergence VI.

 $\rightarrow$  Question Can we find some limitations to this approach?

## A key example

#### Log-normal distribution of the relative weights

Let  $\sigma > 0$ ,  $S_1, \ldots, S_N$  be i.i.d. normal r.v and assume that the distribution of the relative weights  $\overline{w}_{\theta,\phi}(z_1), \ldots, \overline{w}_{\theta,\phi}(z_N)$  is log-normal of the form

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0,1)$ ,

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right)$$

with

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) = -\frac{\alpha\sigma^2 d}{2} \quad \text{and} \quad \gamma_\alpha^2 = \frac{\exp\left[(1-\alpha)^2\sigma^2 d\right] - 1}{1-\alpha}.$$

 $\rightarrow$  Our theorem may not capture what is happening in **high dimensions** i.e. we **may never use** N **large enough** in high-dimensional settings for the asymptotic regime to kick in

 $\rightarrow$  Question Analysis as both d and N go to infinity?  $\Delta_{N,d}^{(\alpha)}(\theta,\phi;x)$
#### Log-normal distribution of the relative weights

Let  $\sigma > 0$ ,  $S_1, \ldots, S_N$  be i.i.d. normal r.v and assume that the distribution of the relative weights  $\overline{w}_{\theta,\phi}(z_1), \ldots, \overline{w}_{\theta,\phi}(z_N)$  is log-normal of the form

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0,1)$ ,

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right)$$

with

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) = -\frac{\alpha\sigma^2 \mathbf{d}}{2} \quad \text{and} \quad \gamma_{\alpha}^2 = \frac{\exp\left[(1-\alpha)^2\sigma^2 d\right] - 1}{1-\alpha}.$$

 $\rightarrow$  Our theorem may not capture what is happening in **high dimensions** i.e. we **may never use** N **large enough** in high-dimensional settings for the asymptotic regime to kick in

#### Log-normal distribution of the relative weights

Let  $\sigma > 0$ ,  $S_1, \ldots, S_N$  be i.i.d. normal r.v and assume that the distribution of the relative weights  $\overline{w}_{\theta,\phi}(z_1), \ldots, \overline{w}_{\theta,\phi}(z_N)$  is log-normal of the form

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0,1)$ ,

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right)$$

with

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) = -\frac{\alpha\sigma^2 \mathbf{d}}{2} \quad \text{and} \quad \gamma_{\alpha}^2 = \frac{\exp\left[(1-\alpha)^2\sigma^2 \mathbf{d}\right] - 1}{1-\alpha}.$$

 $\rightarrow$  Our theorem may not capture what is happening in **high dimensions** i.e. we **may never use** N **large enough** in high-dimensional settings for the asymptotic regime to kick in

#### Log-normal distribution of the relative weights

Let  $\sigma > 0$ ,  $S_1, \ldots, S_N$  be i.i.d. normal r.v and assume that the distribution of the relative weights  $\overline{w}_{\theta,\phi}(z_1), \ldots, \overline{w}_{\theta,\phi}(z_N)$  is log-normal of the form

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0,1)$ ,

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right)$$

with

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) = -\frac{\alpha\sigma^2 \mathbf{d}}{2} \quad \text{and} \quad \gamma_{\alpha}^2 = \frac{\exp\left[(1-\alpha)^2\sigma^2 \mathbf{d}\right] - 1}{1-\alpha}.$$

 $\rightarrow$  Our theorem may not capture what is happening in **high dimensions** i.e. we may never use N large enough in high-dimensional settings for the asymptotic regime to kick in

#### Log-normal distribution of the relative weights

Let  $\sigma > 0$ ,  $S_1, \ldots, S_N$  be i.i.d. normal r.v and assume that the distribution of the relative weights  $\overline{w}_{\theta,\phi}(z_1), \ldots, \overline{w}_{\theta,\phi}(z_N)$  is log-normal of the form

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0,1)$ ,

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right)$$

with

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) = -\frac{\alpha\sigma^2 \mathbf{d}}{2} \quad \text{and} \quad \gamma_{\alpha}^2 = \frac{\exp\left[(1-\alpha)^2\sigma^2 \mathbf{d}\right] - 1}{1-\alpha}.$$

 $\rightarrow$  Our theorem may not capture what is happening in **high dimensions** i.e. we may never use N large enough in high-dimensional settings for the asymptotic regime to kick in

#### Log-normal distribution of the relative weights

Let  $\sigma > 0$ ,  $S_1, \ldots, S_N$  be i.i.d. normal r.v and assume that the distribution of the relative weights  $\overline{w}_{\theta,\phi}(z_1), \ldots, \overline{w}_{\theta,\phi}(z_N)$  is log-normal of the form

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0,1)$ ,

$$\Delta_N^{(\alpha)}(\theta,\phi;x) = \mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) - \frac{\gamma_\alpha^2}{2N} + o\left(\frac{1}{N}\right)$$

with

$$\mathcal{L}^{(\alpha)}(\theta,\phi;x) - \ell(\theta;x) = -\frac{\alpha\sigma^2 \mathbf{d}}{2} \quad \text{and} \quad \gamma_{\alpha}^2 = \frac{\exp\left[(1-\alpha)^2\sigma^2 \mathbf{d}\right] - 1}{1-\alpha}.$$

 $\rightarrow$  Our theorem may not capture what is happening in **high dimensions** i.e. we may never use N large enough in high-dimensional settings for the asymptotic regime to kick in

# Outline

## 1 Introduction

## 2 The VR-IWAE bound

### **3** Theoretical study of the VR-IWAE bound

- Overview First study Second study
- 4 Numerical experiments

## **5** Conclusion

$$N, d 
ightarrow \infty$$
 with either  $rac{\log N}{d} 
ightarrow 0$  or  $rac{\log N}{d^{1/3}} 
ightarrow 0$ 

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N.$$

ightarrow Theoretical study in two steps :

- $\textbf{O} \text{ Log-normal case}: d, N \to \infty \text{ with } \frac{\log N}{d} \to 0$
- **②** Approximate log-normal case :  $d, N \to \infty$  with  $\frac{\log N}{d^{1/3}} \to 0$

$$N, d 
ightarrow \infty$$
 with either  $rac{\log N}{d} 
ightarrow 0$  or  $rac{\log N}{d^{1/3}} 
ightarrow 0$ 

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N.$$

 $\rightarrow$  Theoretical study in two steps :

- $\textbf{0} \text{ Log-normal case}: \ d, N \to \infty \text{ with } \frac{\log N}{d} \to 0$
- **2** Approximate log-normal case :  $d, N \to \infty$  with  $\frac{\log N}{d^{1/3}} \to 0$

$$N, d 
ightarrow \infty$$
 with either  $rac{\log N}{d} 
ightarrow 0$  or  $rac{\log N}{d^{1/3}} 
ightarrow 0$ 

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N.$$

 $\rightarrow$  Theoretical study in two steps :

 $\textbf{O} \text{ Log-normal case}: \ d, N \to \infty \text{ with } \frac{\log N}{d} \to 0$ 

**②** Approximate log-normal case :  $d, N \to \infty$  with  $\frac{\log N}{d^{1/3}} \to 0$ 

$$N, d o \infty$$
 with either  $rac{\log N}{d} o 0$  or  $rac{\log N}{d^{1/3}} o 0$ 

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N.$$

 $\rightarrow$  Theoretical study in two steps :

- $\textbf{O} \text{ Log-normal case }: \ d, N \to \infty \text{ with } \frac{\log N}{d} \to 0$
- **2** Approximate log-normal case :  $d, N \to \infty$  with  $\frac{\log N}{d^{1/3}} \to 0$

#### Theorem

Let  $S_1,\ldots,S_N$  be i.i.d. normal random variables. Further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0,1)$ , we have

$$\lim_{\substack{N,d\to\infty\\ \log N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left( 1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right) \right) = 0.$$

ightarrow Previously, we had

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

- While increasing N decreases the variational gap for N large enough, it does so by a factor which is **negligible** before the term  $-d\sigma^2/2$
- $\bullet\,$  This time, the term  $-d\sigma^2/2$  does not depend on  $\alpha$
- ightarrow Weight collapse phenomenon : for all  $lpha \in [0,1)$ ,

 $\Delta_{N,d}^{(\alpha)}(\theta,\phi;x)\approx \mathrm{ELBO}(\theta,\phi;x)-\ell(\theta;x), \quad \text{as } N,d\rightarrow\infty \text{ with } \frac{\log N}{d}\rightarrow 0.$ 

#### Theorem

Let  $S_1,\ldots,S_N$  be i.i.d. normal random variables. Further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left( 1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right) \right) = 0.$$

 $\rightarrow$  Previously, we had

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

- While increasing N decreases the variational gap for N large enough, it does so by a factor which is **negligible** before the term  $-d\sigma^2/2$
- $\bullet\,$  This time, the term  $-d\sigma^2/2$  does not depend on  $\alpha$
- ightarrow Weight collapse phenomenon : for all  $lpha \in [0,1)$ ,

 $\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) \approx \text{ELBO}(\theta,\phi;x) - \ell(\theta;x), \quad \text{as } N, d \to \infty \text{ with } \frac{\log N}{d} \to 0.$ 

#### Theorem

Let  $S_1,\ldots,S_N$  be i.i.d. normal random variables. Further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left( 1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right) \right) = 0.$$

ightarrow Previously, we had

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

- While increasing N decreases the variational gap for N large enough, it does so by a factor which is **negligible** before the term  $-d\sigma^2/2$
- This time, the term  $-d\sigma^2/2$  does not depend on lpha
- $\rightarrow$  Weight collapse phenomenon : for all  $\alpha \in [0, 1)$ ,

 $\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) \approx \text{ELBO}(\theta,\phi;x) - \ell(\theta;x), \quad \text{as } N, d \to \infty \text{ with } \frac{\log N}{d} \to 0.$ 

#### Theorem

Let  $S_1,\ldots,S_N$  be i.i.d. normal random variables. Further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left( 1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right) \right) = 0.$$

 $\rightarrow$  Previously, we had

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

- While increasing N decreases the variational gap for N large enough, it does so by a factor which is **negligible** before the term  $-d\sigma^2/2$
- This time, the term  $-d\sigma^2/2~{\rm does~not}$  depend on  $\alpha$
- $\rightarrow$  Weight collapse phenomenon : for all  $\alpha \in [0, 1)$ ,

 $\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) \approx \text{ELBO}(\theta,\phi;x) - \ell(\theta;x), \quad \text{as } N, d \to \infty \text{ with } \frac{\log N}{d} \to 0.$ 

#### Theorem

Let  $S_1,\ldots,S_N$  be i.i.d. normal random variables. Further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad i = 1 \dots N.$$

Then, for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left( 1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right) \right) = 0.$$

ightarrow Previously, we had

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

- While increasing N decreases the variational gap for N large enough, it does so by a factor which is **negligible** before the term  $-d\sigma^2/2$
- This time, the term  $-d\sigma^2/2~{\rm does~not}$  depend on  $\alpha$
- $\rightarrow$  Weight collapse phenomenon : for all  $\alpha \in [0,1),$

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) \approx \mathrm{ELBO}(\theta,\phi;x) - \ell(\theta;x), \quad \text{as } N,d \to \infty \text{ with } \frac{\log N}{d} \to 0.$$

#### Gaussian example

Set  $p_{\theta}(z|x) = \mathcal{N}(z; \theta, \mathbf{I}_d)$  and  $q_{\phi}(z) = \mathcal{N}(z; \phi, \mathbf{I}_d)$ , with  $\theta = 0 \cdot \mathbf{u}_d$  and  $\phi = \mathbf{u}_d$ , where  $\mathbf{u}_d$  is the *d*-dimensional vector whose coordinates are all equal to 1. Then

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N$$

with  $\sigma = 1$ .

• Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

• Asymptotic result 2

$$\lim_{\substack{N,d\to\infty\\\log N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left(1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

 $\rightarrow$  Weight collapse phenomenon might occur even for simple examples!

#### Gaussian example

Set  $p_{\theta}(z|x) = \mathcal{N}(z; \theta, \mathbf{I}_d)$  and  $q_{\phi}(z) = \mathcal{N}(z; \phi, \mathbf{I}_d)$ , with  $\theta = 0 \cdot \mathbf{u}_d$  and  $\phi = \mathbf{u}_d$ , where  $\mathbf{u}_d$  is the *d*-dimensional vector whose coordinates are all equal to 1. Then

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N$$

with  $\sigma = 1$ .

Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

• Asymptotic result 2

$$\lim_{\substack{N,d\to\infty\\ {\rm og}\,N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left(1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

ightarrow Weight collapse phenomenon might occur even for simple examples!

#### Gaussian example

Set  $p_{\theta}(z|x) = \mathcal{N}(z; \theta, \mathbf{I}_d)$  and  $q_{\phi}(z) = \mathcal{N}(z; \phi, \mathbf{I}_d)$ , with  $\theta = 0 \cdot \mathbf{u}_d$  and  $\phi = \mathbf{u}_d$ , where  $\mathbf{u}_d$  is the *d*-dimensional vector whose coordinates are all equal to 1. Then

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N$$

with  $\sigma = 1$ .

Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

• Asymptotic result 2

$$\lim_{\substack{N,d \to \infty \\ \log N/d \to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left( 1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha} \frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right) \right) = 0.$$

ightarrow Weight collapse phenomenon might occur even for simple examples!

#### Gaussian example

Set  $p_{\theta}(z|x) = \mathcal{N}(z; \theta, \mathbf{I}_d)$  and  $q_{\phi}(z) = \mathcal{N}(z; \phi, \mathbf{I}_d)$ , with  $\theta = 0 \cdot \mathbf{u}_d$  and  $\phi = \mathbf{u}_d$ , where  $\mathbf{u}_d$  is the *d*-dimensional vector whose coordinates are all equal to 1. Then

$$\log \overline{w}_{\theta,\phi}(z_i) = -\frac{\sigma^2 d}{2} - \sigma \sqrt{d} S_i, \quad S_i \sim \mathcal{N}(0,1), \quad i = 1 \dots N$$

with  $\sigma = 1$ .

Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = -\alpha \cdot \frac{\sigma^2 d}{2} - \frac{\exp\left[(1-\alpha)^2 \sigma^2 d\right] - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

• Asymptotic result 2

$$\lim_{\substack{N,d\to\infty\\\log N/d\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + \frac{d\sigma^2}{2} \left(1 - 2\sqrt{\frac{2\log N}{d\sigma^2}} + \frac{1}{1-\alpha}\frac{2\log N}{d\sigma^2} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

 $\rightarrow$  Weight collapse phenomenon might occur even for simple examples!

# Empirical verification



# Empirical verification (cont'd)



















Challenges and Opportunities in Scalable  $\alpha$ -divergence Variational Inference 26 / 37

Kamélia Daudel (University of Oxford)

(A1) For all  $i = 1 \dots N$ ,

- $\xi_{i,1}, \ldots, \xi_{i,d}$  are i.i.d. random variables which are absolutely continuous with respect to the Lebesgue measure and satisfy  $\mathbb{E}(\xi_{i,1}) = 0$  and  $\mathbb{V}(\xi_{i,1}) = \sigma^2 < \infty$ .
- **2** There exists K > 0 such that:

$$|\mathbb{E}(\xi_{i,1}^k)| \le k! K^{k-2} \sigma^2, \quad k \ge 3.$$

 $\rightarrow$  Let  $S_1, \ldots, S_N$  be such that :

$$S_i = \frac{1}{\sigma\sqrt{d}} \sum_{j=1}^d \xi_{i,j}, \quad i = 1\dots N.$$

#### Theorem

Assume (A1). Set  $a := \log \mathbb{E}(\exp(-\xi_{1,1}))$  and further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$

Then, a > 0 and for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0$$

$$\rightarrow \text{NB}: -da = -\log \mathbb{E}(\exp(-\sigma\sqrt{dS_1}))$$

(A1) For all  $i = 1 \dots N$ ,

- $\xi_{i,1}, \ldots, \xi_{i,d}$  are i.i.d. random variables which are absolutely continuous with respect to the Lebesgue measure and satisfy  $\mathbb{E}(\xi_{i,1}) = 0$  and  $\mathbb{V}(\xi_{i,1}) = \sigma^2 < \infty$ .
- **2** There exists K > 0 such that:

$$|\mathbb{E}(\xi_{i,1}^k)| \le k! K^{k-2} \sigma^2, \quad k \ge 3.$$

 $\rightarrow$  Let  $S_1, \ldots, S_N$  be such that :

$$S_i = \frac{1}{\sigma\sqrt{d}} \sum_{j=1}^d \xi_{i,j}, \quad i = 1\dots N.$$

#### Theorem

Assume (A1). Set  $a := \log \mathbb{E}(\exp(-\xi_{1,1}))$  and further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$

Then, a > 0 and for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0$$

 $\rightarrow \mathsf{NB} : -da = -\log \mathbb{E}(\exp(-\sigma \sqrt{dS_1}))$ 

(A1) For all  $i = 1 \dots N$ ,

- $\xi_{i,1}, \ldots, \xi_{i,d}$  are i.i.d. random variables which are absolutely continuous with respect to the Lebesgue measure and satisfy  $\mathbb{E}(\xi_{i,1}) = 0$  and  $\mathbb{V}(\xi_{i,1}) = \sigma^2 < \infty$ .
- **2** There exists K > 0 such that:

$$|\mathbb{E}(\xi_{i,1}^k)| \le k! K^{k-2} \sigma^2, \quad k \ge 3.$$

 $\rightarrow$  Let  $S_1, \ldots, S_N$  be such that :

$$S_i = \frac{1}{\sigma\sqrt{d}} \sum_{j=1}^d \xi_{i,j}, \quad i = 1\dots N.$$

#### Theorem

Assume (A1). Set  $a := \log \mathbb{E}(\exp(-\xi_{1,1}))$  and further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$

Then, a > 0 and for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\ \log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

 $\rightarrow \text{NB}: -da = -\log \mathbb{E}(\exp(-\sigma \sqrt{dS_1}))$ 

(A1) For all  $i = 1 \dots N$ ,

- $\xi_{i,1}, \ldots, \xi_{i,d}$  are i.i.d. random variables which are absolutely continuous with respect to the Lebesgue measure and satisfy  $\mathbb{E}(\xi_{i,1}) = 0$  and  $\mathbb{V}(\xi_{i,1}) = \sigma^2 < \infty$ .
- **2** There exists K > 0 such that:

$$|\mathbb{E}(\xi_{i,1}^k)| \le k! K^{k-2} \sigma^2, \quad k \ge 3.$$

 $\rightarrow$  Let  $S_1, \ldots, S_N$  be such that :

$$S_i = \frac{1}{\sigma\sqrt{d}} \sum_{j=1}^d \xi_{i,j}, \quad i = 1\dots N.$$

#### Theorem

Assume (A1). Set  $a := \log \mathbb{E}(\exp(-\xi_{1,1}))$  and further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$

Then, a > 0 and for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

$$\rightarrow \mathsf{NB} : -da = -\log \mathbb{E}(\exp(-\sigma \sqrt{d}S_1))$$

# Main result in the approximate log-normal case (cont'd) $\rightarrow$ Let $S_1, \ldots, S_N$ be such that :

$$S_i = \frac{1}{\sigma\sqrt{d}} \sum_{j=1}^d \xi_{i,j}, \quad i = 1\dots N.$$

#### Theorem

Assume (A1). Set  $a := \log \mathbb{E}(\exp(-\xi_{1,1}))$  and further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{dS_i}, \quad i = 1 \dots N.$$

Then, a > 0 and for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

ightarrow Weight collapse phenomenon : for all  $lpha \in [0,1)$ ,

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x)\approx \mathrm{ELBO}(\theta,\phi;x)-\ell(\theta;x), \quad \text{as } N,d\rightarrow\infty \text{ with } \tfrac{\log N}{d^{1/3}}\rightarrow 0.$$

The condition that N should grow at least exponentially with d has been replaced by the less restrictive yet still stringent condition that N should grow at least sub-exponentially with  $d^{1/3}$ .

ightarrow NB : no dependency in lpha left in the asymptotic regime

# Main result in the approximate log-normal case (cont'd) $\rightarrow$ Let $S_1, \ldots, S_N$ be such that :

$$S_i = \frac{1}{\sigma\sqrt{d}} \sum_{j=1}^d \xi_{i,j}, \quad i = 1\dots N.$$

#### Theorem

Assume (A1). Set  $a := \log \mathbb{E}(\exp(-\xi_{1,1}))$  and further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$

Then, a > 0 and for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

 $\rightarrow$  Weight collapse phenomenon : for all  $\alpha \in [0, 1)$ ,

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x)\approx \mathrm{ELBO}(\theta,\phi;x)-\ell(\theta;x), \quad \text{as } N,d\rightarrow\infty \text{ with } \frac{\log N}{d^{1/3}}\rightarrow 0.$$

The condition that N should grow at least exponentially with d has been replaced by the less restrictive yet still stringent condition that N should grow at least sub-exponentially with  $d^{1/3}$ .

ightarrow NB : no dependency in lpha left in the asymptotic regime

# Main result in the approximate log-normal case (cont'd) $\rightarrow$ Let $S_1, \ldots, S_N$ be such that :

$$S_i = \frac{1}{\sigma\sqrt{d}} \sum_{j=1}^d \xi_{i,j}, \quad i = 1\dots N.$$

#### Theorem

Assume (A1). Set  $a := \log \mathbb{E}(\exp(-\xi_{1,1}))$  and further assume that

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{dS_i}, \quad i = 1 \dots N.$$

Then, a > 0 and for all  $\alpha \in [0, 1)$ , we have

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

 $\rightarrow$  Weight collapse phenomenon : for all  $\alpha \in [0, 1)$ ,

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x)\approx \mathrm{ELBO}(\theta,\phi;x)-\ell(\theta;x),\quad \text{as }N,d\rightarrow\infty \text{ with } \frac{\log N}{d^{1/3}}\rightarrow 0.$$

The condition that N should grow at least exponentially with d has been replaced by the less restrictive yet still stringent condition that N should grow at least sub-exponentially with  $d^{1/3}$ .

 $\rightarrow$  NB : no dependency in  $\alpha$  left in the asymptotic regime

#### Linear Gaussian example (Rainforth et al., 2018)

Set  $p_{\theta}(z) = \mathcal{N}(z; \theta, \mathbf{I}_d)$ ,  $p_{\theta}(x|z) = \mathcal{N}(x; z, \mathbf{I}_d)$  with  $\theta \in \mathbb{R}^d$ , and  $q_{\phi}(z|x) = \mathcal{N}(z; Ax + b, 2/3 \ \mathbf{I}_d)$  with  $A = \operatorname{diag}(\tilde{a})$  and  $\phi = (\tilde{a}, b) \in \mathbb{R}^d \times \mathbb{R}^d$ . Then, we can write

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$
  
with  $\sigma^2 = \frac{1}{18} + \frac{8}{3}\lambda^2$  and  $a = \lambda^2 + \frac{1}{6} + \frac{1}{2}\log(3/4)$ , where  $\lambda = \frac{\left\|\frac{\theta+x}{2} - Ax - b\right\|}{\sqrt{d}}$ 

$$\rightarrow \mathsf{Set}\left(\theta,\phi\right) = \left(\theta^{\star},\phi^{\star}\right)\left[\theta^{\star} = T^{-1}\sum_{t=1}^{T} x_{t}, \phi^{\star} = \left(a^{\star},b^{\star}\right) \text{ with } a^{\star} = \frac{1}{2}u_{d}, \ b^{\star} = \frac{\theta^{\star}}{2}\right]$$

• Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = \frac{d}{2} \left[ \log\left(\frac{4}{3}\right) + \frac{1}{1-\alpha} \log\left(\frac{3}{4-\alpha}\right) \right] - \frac{(4-\alpha)^d (15-6\alpha)^{-\frac{d}{2}} - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

• Asymptotic result 2

$$\lim_{\substack{N,d \to \infty \\ \log N/d^{1/3} \to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0$$

#### ightarrow The choice of the variational approximation $q_{\phi}$ matters a lot!

#### Linear Gaussian example (Rainforth et al., 2018)

Set  $p_{\theta}(z) = \mathcal{N}(z; \theta, \mathbf{I}_d)$ ,  $p_{\theta}(x|z) = \mathcal{N}(x; z, \mathbf{I}_d)$  with  $\theta \in \mathbb{R}^d$ , and  $q_{\phi}(z|x) = \mathcal{N}(z; Ax + b, 2/3 \ \mathbf{I}_d)$  with  $A = \operatorname{diag}(\tilde{a})$  and  $\phi = (\tilde{a}, b) \in \mathbb{R}^d \times \mathbb{R}^d$ . Then, we can write

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$
  
with  $\sigma^2 = \frac{1}{18} + \frac{8}{3}\lambda^2$  and  $a = \lambda^2 + \frac{1}{6} + \frac{1}{2}\log(3/4)$ , where  $\lambda = \frac{\left\|\frac{\theta+x}{2} - Ax - b\right\|}{\sqrt{d}}$ 

$$\rightarrow \mathsf{Set} \left( \theta, \phi \right) = \left( \theta^{\star}, \phi^{\star} \right) \left[ \theta^{\star} = T^{-1} \sum_{t=1}^{T} x_t, \, \phi^{\star} = \left( a^{\star}, b^{\star} \right) \text{ with } a^{\star} = \frac{1}{2} \boldsymbol{u}_d, \, b^{\star} = \frac{\theta^{\star}}{2} \right]$$

• Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = \frac{d}{2} \left[ \log\left(\frac{4}{3}\right) + \frac{1}{1-\alpha} \log\left(\frac{3}{4-\alpha}\right) \right] - \frac{(4-\alpha)^d (15-6\alpha)^{-\frac{d}{2}} - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

• Asymptotic result 2

$$\lim_{\substack{N,d\to\infty\\ \log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0$$

#### ightarrow The choice of the variational approximation $q_{\phi}$ matters a lot!

#### Linear Gaussian example (Rainforth et al., 2018)

Set  $p_{\theta}(z) = \mathcal{N}(z; \theta, \mathbf{I}_d)$ ,  $p_{\theta}(x|z) = \mathcal{N}(x; z, \mathbf{I}_d)$  with  $\theta \in \mathbb{R}^d$ , and  $q_{\phi}(z|x) = \mathcal{N}(z; Ax + b, 2/3 \ \mathbf{I}_d)$  with  $A = \operatorname{diag}(\tilde{a})$  and  $\phi = (\tilde{a}, b) \in \mathbb{R}^d \times \mathbb{R}^d$ . Then, we can write

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$
  
with  $\sigma^2 = \frac{1}{18} + \frac{8}{3}\lambda^2$  and  $a = \lambda^2 + \frac{1}{6} + \frac{1}{2}\log(3/4)$ , where  $\lambda = \frac{\left\|\frac{\theta+x}{2} - Ax - b\right\|}{\sqrt{d}}$ 

$$\rightarrow \mathsf{Set} \left( \theta, \phi \right) = \left( \theta^{\star}, \phi^{\star} \right) \left[ \theta^{\star} = T^{-1} \sum_{t=1}^{T} x_t, \, \phi^{\star} = \left( a^{\star}, b^{\star} \right) \text{ with } a^{\star} = \frac{1}{2} \boldsymbol{u}_d, \, b^{\star} = \frac{\theta^{\star}}{2} \right]$$

• Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = \frac{d}{2} \left[ \log\left(\frac{4}{3}\right) + \frac{1}{1-\alpha} \log\left(\frac{3}{4-\alpha}\right) \right] - \frac{(4-\alpha)^d (15-6\alpha)^{-\frac{d}{2}} - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

• Asymptotic result 2

$$\lim_{\substack{N,d\to\infty\\ \log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0$$

 $\rightarrow$  The choice of the variational approximation  $q_{\phi}$  matters a lot!

#### Linear Gaussian example (Rainforth et al., 2018)

Set  $p_{\theta}(z) = \mathcal{N}(z; \theta, \mathbf{I}_d)$ ,  $p_{\theta}(x|z) = \mathcal{N}(x; z, \mathbf{I}_d)$  with  $\theta \in \mathbb{R}^d$ , and  $q_{\phi}(z|x) = \mathcal{N}(z; Ax + b, 2/3 \ \mathbf{I}_d)$  with  $A = \operatorname{diag}(\tilde{a})$  and  $\phi = (\tilde{a}, b) \in \mathbb{R}^d \times \mathbb{R}^d$ . Then, we can write

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$
  
with  $\sigma^2 = \frac{1}{18} + \frac{8}{3}\lambda^2$  and  $a = \lambda^2 + \frac{1}{6} + \frac{1}{2}\log(3/4)$ , where  $\lambda = \frac{\left\|\frac{\theta+x}{2} - Ax - b\right\|}{\sqrt{d}}$ 

$$\rightarrow \mathsf{Set} (\theta, \phi) = (\theta^{\star}, \phi^{\star}) \left[ \theta^{\star} = T^{-1} \sum_{t=1}^{T} x_t, \, \phi^{\star} = (a^{\star}, b^{\star}) \text{ with } a^{\star} = \frac{1}{2} \boldsymbol{u}_d, \, b^{\star} = \frac{\theta^{\star}}{2} \right]$$

• Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = \frac{d}{2} \left[ \log\left(\frac{4}{3}\right) + \frac{1}{1-\alpha} \log\left(\frac{3}{4-\alpha}\right) \right] - \frac{(4-\alpha)^d (15-6\alpha)^{-\frac{d}{2}} - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

Asymptotic result 2

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}}\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

 $\rightarrow$  The choice of the variational approximation  $q_{\phi}$  matters a lot!

#### Linear Gaussian example (Rainforth et al., 2018)

Set  $p_{\theta}(z) = \mathcal{N}(z; \theta, \mathbf{I}_d)$ ,  $p_{\theta}(x|z) = \mathcal{N}(x; z, \mathbf{I}_d)$  with  $\theta \in \mathbb{R}^d$ , and  $q_{\phi}(z|x) = \mathcal{N}(z; Ax + b, 2/3 \ \mathbf{I}_d)$  with  $A = \operatorname{diag}(\tilde{a})$  and  $\phi = (\tilde{a}, b) \in \mathbb{R}^d \times \mathbb{R}^d$ . Then, we can write

$$\log \overline{w}_{\theta,\phi}(z_i) = -da - \sigma \sqrt{d}S_i, \quad i = 1 \dots N.$$
  
with  $\sigma^2 = \frac{1}{18} + \frac{8}{3}\lambda^2$  and  $a = \lambda^2 + \frac{1}{6} + \frac{1}{2}\log(3/4)$ , where  $\lambda = \frac{\left\|\frac{\theta+x}{2} - Ax - b\right\|}{\sqrt{d}}$ 

$$\rightarrow \mathsf{Set}\;(\theta,\phi) = (\theta^{\star},\phi^{\star})\;[\theta^{\star} = T^{-1}\sum_{t=1}^{T} x_t,\,\phi^{\star} = (a^{\star},b^{\star})\;\mathsf{with}\;a^{\star} = \frac{1}{2}\boldsymbol{u}_d,\;b^{\star} = \frac{\theta^{\star}}{2}]$$

• Asymptotic result 1

$$\Delta_{N,d}^{(\alpha)}(\theta,\phi;x) = \frac{d}{2} \left[ \log\left(\frac{4}{3}\right) + \frac{1}{1-\alpha} \log\left(\frac{3}{4-\alpha}\right) \right] - \frac{(4-\alpha)^d (15-6\alpha)^{-\frac{d}{2}} - 1}{2(1-\alpha)N} + o\left(\frac{1}{N}\right)$$

Asymptotic result 2

$$\lim_{\substack{N,d\to\infty\\\log N/d^{1/3}\to 0}} \Delta_{N,d}^{(\alpha)}(\theta,\phi;x) + da\left(1 - \frac{\sigma}{a}\sqrt{\frac{2\log N}{d}} + O\left(\frac{\log\log N}{\sqrt{d\log N}}\right)\right) = 0.$$

ightarrow The choice of the variational approximation  $q_{\phi}$  matters a lot!

# Outline

## 1 Introduction

- 2 The VR-IWAE bound
- 3 Theoretical study of the VR-IWAE bound
- **4** Numerical experiments

### **5** Conclusion

# Linear Gaussian example



Kamélia Daudel (University of Oxford) · Challenges and Op

Challenges and Opportunities in Scalable  $\alpha$ -divergence Variational Inference 31 / 37

# Linear Gaussian example (cont'd)


# Variational auto-encoder on MNIST



# Variational auto-encoder on MNIST (cont'd)



Kamélia Daudel (University of Oxford) · Challenges and Opportunities in Scalable  $\alpha$ -divergence Variational Inference 34 / 37

#### Variational auto-encoder on MNIST (cont'd - 2)



# Outline

#### 1 Introduction

- 2 The VR-IWAE bound
- 3 Theoretical study of the VR-IWAE bound
- **4** Numerical experiments



Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)
- <sup>(2)</sup> We provided two complementary analyses of the VR-IWAR bound
  - Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
  - Encompass the case of the IWAE bound
- 3 Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)
- <sup>(2)</sup> We provided two complementary analyses of the VR-IWAR bound
  - Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
  - Encompass the case of the IWAE bound
- 3 Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)
- We provided two complementary analyses of the VR-IWAR bound
  - Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
  - Encompass the case of the IWAE bound
- **③** Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)
- <sup>(2)</sup> We provided two complementary analyses of the VR-IWAR bound
  - Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
  - Encompass the case of the IWAE bound
- **③** Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

#### We provided two complementary analyses of the VR-IWAR bound

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound

#### 3 Empirical verification of our theoretical results

- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound
- 3 Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound
- 3 Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

**2** We provided two complementary analyses of the VR-IWAR bound

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound

#### **3** Empirical verification of our theoretical results

#### $\rightarrow$ Further work:

- Does the weight collapse behavior apply beyond the cases highlighted here?
- How does the weight collapse affect the gradient descent procedures?
- Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

**2** We provided two complementary analyses of the VR-IWAR bound

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound

#### **3** Empirical verification of our theoretical results

#### $\rightarrow$ Further work:

- Does the weight collapse behavior apply beyond the cases highlighted here?
- How does the weight collapse affect the gradient descent procedures?
- Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound
- **3** Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound
- **3** Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?

Daudel, Benton, Shi and Doucet (2022). Alpha-divergence Variational Inference Meets Importance Weighted Auto-Encoders: Methodology and Asymptotics.

• We formalized and motivated the VR-IWAE bound

- Theoretically-sound extension of the IWAE bound ( $\alpha = 0$ )
- Provides theoretical guarantees behind various VR bound-based schemes proposed in the  $\alpha$ -Divergence VI community
- Enjoys other additional desirable properties of this bound (SNR, doubly-reparameterized)

- Shed light on the conditions behind the success or failure of the VR-IWAE bound methodology
- Encompass the case of the IWAE bound
- **3** Empirical verification of our theoretical results
- $\rightarrow$  Further work:
  - Does the weight collapse behavior apply beyond the cases highlighted here?
  - How does the weight collapse affect the gradient descent procedures?
  - Can we further use the fact that the VR-IWAE bound extends the IWAE bound?