

Motivation

Compute / sample from the **posterior density** of the latent variables y given the data \mathcal{D} :

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})}.$$

Problem: the marginal likelihood $p(\mathcal{D})$ is **untractable**.

→ **Variational Inference** methods see this as an **optimisation problem** over the variational family $\{y \mapsto q_\theta(y) : \theta \in \mathcal{T}\}$.

Let us now consider a broader approximating family

$$\left\{ y \mapsto \int_{\mathcal{T}} \mu(d\theta) q_\theta(y) : \mu \in \mathcal{M} \right\},$$

where \mathcal{M} is a subset of $\mathcal{M}_1(\mathcal{T})$, the set of probability measures on $(\mathcal{T}, \mathcal{T})$.

Question: Can we define an **iterative scheme** which diminishes a given objective function at each step ? → **Yes** : the f -EI(ϕ) Algorithm !

The f -EI(ϕ) Algorithm

General Optimisation Problem: f convex over $(0, \infty)$, $f(1) = 0$

$$\operatorname{arginf}_{\mu \in \mathcal{M}} \Psi^{(f)}(\mu) \quad \text{where} \quad \Psi^{(f)}(\mu) = \int_{\mathcal{Y}} f\left(\frac{\mu q(y)}{p(y)}\right) p(y) \nu(dy).$$

→ the mapping $\mu \mapsto \Psi^{(f)}(\mu)$ is **convex**.

→ includes **f -Divergence** posterior density approximation.

Let $\phi \in \mathbb{R}^*$, $\mu \in \mathcal{M}_1(\mathcal{T})$ and let the sequence $(\mu_n)_{n \in \mathbb{N}}$ be defined by

$$\begin{cases} \mu_0 = \mu, \\ \mu_{n+1} = \mathcal{I}^\phi(\mu_n), \end{cases} \quad n \in \mathbb{N}. \quad (1)$$

where for all $\zeta \in \mathcal{M}_1(\mathcal{T})$,

$$1. \text{ Expectation step : } b_\zeta(\theta) = \int_{\mathcal{Y}} q(\theta, y) f'\left(\frac{\zeta q(y)}{p(y)}\right) \nu(dy),$$

$$2. \text{ Iteration step : } \mathcal{I}^\phi(\zeta)(d\theta) = \frac{\zeta(d\theta) \cdot |b_\zeta(\theta)|^\phi}{\zeta(|b_\zeta|^\phi)}.$$

(A1) For all $(\theta, y) \in \mathcal{T} \times \mathcal{Y}$, $q(\theta, y) > 0$, $p(y) > 0$ and $\int_{\mathcal{Y}} p(y) \nu(dy) < \infty$.

(A2) $f : (0, \infty) \rightarrow \mathbb{R}$ is monotonous, strictly convex and continuously differentiable, and $f(1) = 0$.

Theoretical Results

Divergence considered		Corresponding range
<i>Reverse KL</i> $f(u) = -\log(u)$		$\phi \in (0, 1]$
<i>α-divergence</i> $f(u) = \frac{1}{\alpha(\alpha-1)}(u^\alpha - 1)$	$\alpha \in (-\infty, -1]$	$\phi \in (0, -1/\alpha]$
	$\alpha \in (-1, 1) \setminus \{0\}$	$\phi \in (0, 1]$
	$\alpha \in (1, \infty)$	$\phi \in (1/(1-\alpha), 0)$

Table 1: Allowed (f, ϕ) in the f -EI(ϕ) algorithm

Theorem 1. Assume (A1). Let (f, ϕ) belong to Table 1. Then (A2) holds. Moreover, let $\mu \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi^{(f)}(\mu) < \infty$. Then the sequence $(\mu_n)_{n \in \mathbb{N}}$ is well-defined and the sequence $(\Psi^{(f)}(\mu_n))_{n \in \mathbb{N}}$ is **non-increasing**.

(A3) \mathcal{T} is a compact metric space, for all $y \in \mathcal{Y}$, $\theta \mapsto q(\theta, y)$ is continuous + uniform boundedness of $\Psi^{(f)}$ and b_μ with respect to μ and θ .

Theorem 2. Assume (A1), (A3) and let (f, ϕ) belong to Table 1. Further assume that there exists $\mu, \bar{\mu} \in \mathcal{M}_1(\mathcal{T})$ such that $\mu_n \Rightarrow \bar{\mu}$ as $n \rightarrow \infty$. Then $\bar{\mu}$ is a fixed point of \mathcal{I}^ϕ and

$$\Psi^{(f)}(\bar{\mu}) = \inf_{\zeta \in \mathcal{M}_1(\mathcal{T})} \Psi^{(f)}(\zeta).$$

Let $Y_1, \dots, Y_K \stackrel{iid}{\sim} \mu q$ and define $b_{\mu, K}(\theta) = \frac{1}{K} \sum_{k=1}^K \frac{q(\theta, Y_k)}{\mu q(Y_k)} f'\left(\frac{\mu q(Y_k)}{p(Y_k)}\right)$.

Theorem 3. Assume (A1). Let (f, ϕ) belong to Table 1. Let $\mu \in \mathcal{M}_1(\mathcal{T})$ be such that $\int_{\mathcal{T}} \mu(d\theta) \mathbb{E}_{\mu q}[\{\frac{q(\theta, Y_1)}{\mu q(Y_1)} | f'\left(\frac{\mu q(Y_1)}{p(Y_1)}\right) | \}^\phi] < \infty$ and $\Psi^{(f)}(\mu) < \infty$. Then, \mathbb{P} - a.s.

$$\lim_{K \rightarrow \infty} \left\| \mathcal{I}_K^\phi(\mu) - \mathcal{I}^\phi(\mu) \right\|_{TV} = 0,$$

where $\mathcal{I}_K^\phi(\mu)(d\theta) = \frac{\mu(d\theta) \cdot |b_{\mu, K}(\theta)|^\phi}{\mu(|b_{\mu, K}|^\phi)}.$

Density Approximation

We can access an **unnormalized** version p^* of the probability density \tilde{p}

$$\tilde{p}(y) = \frac{p^*(y)}{Z}, \quad \text{where} \quad Z := \int_{\mathcal{Y}} p^*(y) \nu(dy).$$

Lemma 1. Assume (A1). For the α -divergence, optimising $D_f(\mu Q || \tilde{\mathbb{P}})$ (with respect to μ) is equivalent to optimising $\Psi^{(f)}(\mu; p)$ with $p = p^*$. Furthermore, for all $\alpha_+ \in (0, 1) \cup (1, +\infty)$ and all $\alpha_- < 0$, we have

$$\forall \mu \in \mathcal{M}_1(\mathcal{T}), \quad \xi^{(\alpha_+)}(\mu q) \leq Z \leq \xi^{(\alpha_-)}(\mu q),$$

$$\text{where } \xi^{(\alpha)}(\tilde{q}) := \left[\int_{\mathcal{Y}} \left(\frac{\tilde{q}(y)}{p^*(y)} \right)^\alpha p^*(y) \nu(dy) \right]^{\frac{1}{1-\alpha}}.$$

Mixture α -Approximate f -EI(ϕ)

Algorithm 1: Mixture α -Approximate f -EI(ϕ)

Input: p^* : unnormalized version of the density \tilde{p} , Q : Markov transition kernel, K : number of samples, $\Theta_J = \{\theta_1, \dots, \theta_J\} \subset \mathcal{T}$: parameter set.

Output: Optimised weights λ .

Set $\lambda = [\frac{1}{J}, \dots, \frac{1}{J}]$.

while the α -bound has not converged **do**

 Sampling step : Draw independently K samples Y_1, \dots, Y_K from $\mu_\lambda q$.

 Expectation step : Compute $\mathbf{A}_\lambda = (a_j)_{1 \leq j \leq J}$ where

$$a_j = \frac{1}{K} \sum_{k=1}^K q(\theta_j, Y_k) \mu_\lambda q(Y_k)^{\alpha-2} p^*(Y_k)^{1-\alpha}$$

 and deduce $\mathbf{B}_\lambda = (\lambda_j a_j^\phi)_{1 \leq j \leq J}$, $b_\lambda = \sum_{j=1}^J \lambda_j a_j^\phi$ and $c_\lambda = \sum_{j=1}^J \lambda_j a_j$.

 Iteration step : Set

$$\xi_K^{(\alpha)}(\mu_\lambda q) \leftarrow c_\lambda^{1/(1-\alpha)}$$

$$\lambda \leftarrow \frac{1}{b_\lambda} \mathbf{B}_\lambda$$

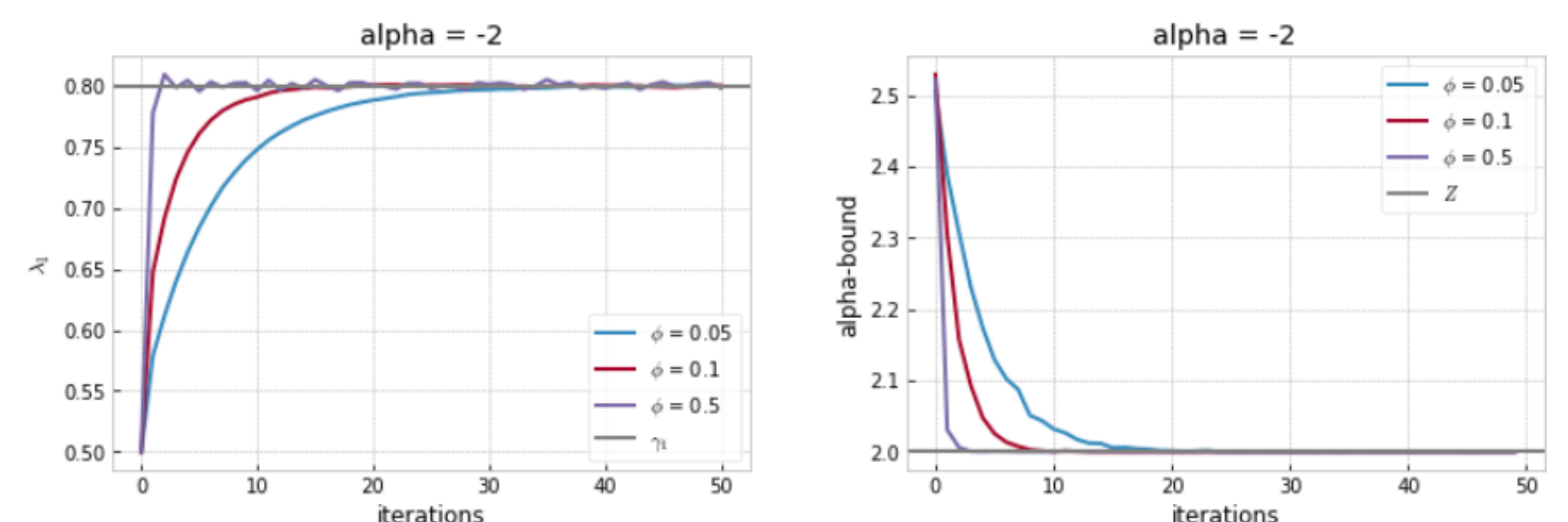
end

→ *Score gradient:* $\tilde{q} \mapsto \mathcal{L}_A^{(\alpha)}(\tilde{q}) := \int_{\mathcal{Y}} \frac{1}{\alpha(\alpha-1)} \left(\frac{\tilde{q}(y)}{p^*(y)} \right)^\alpha p^*(y) \nu(dy),$

$$\nabla_\lambda \mathcal{L}_A^{(\alpha)}(\mu_\lambda q) = (b_{\mu_\lambda}(\theta_j))_{1 \leq j \leq J}.$$

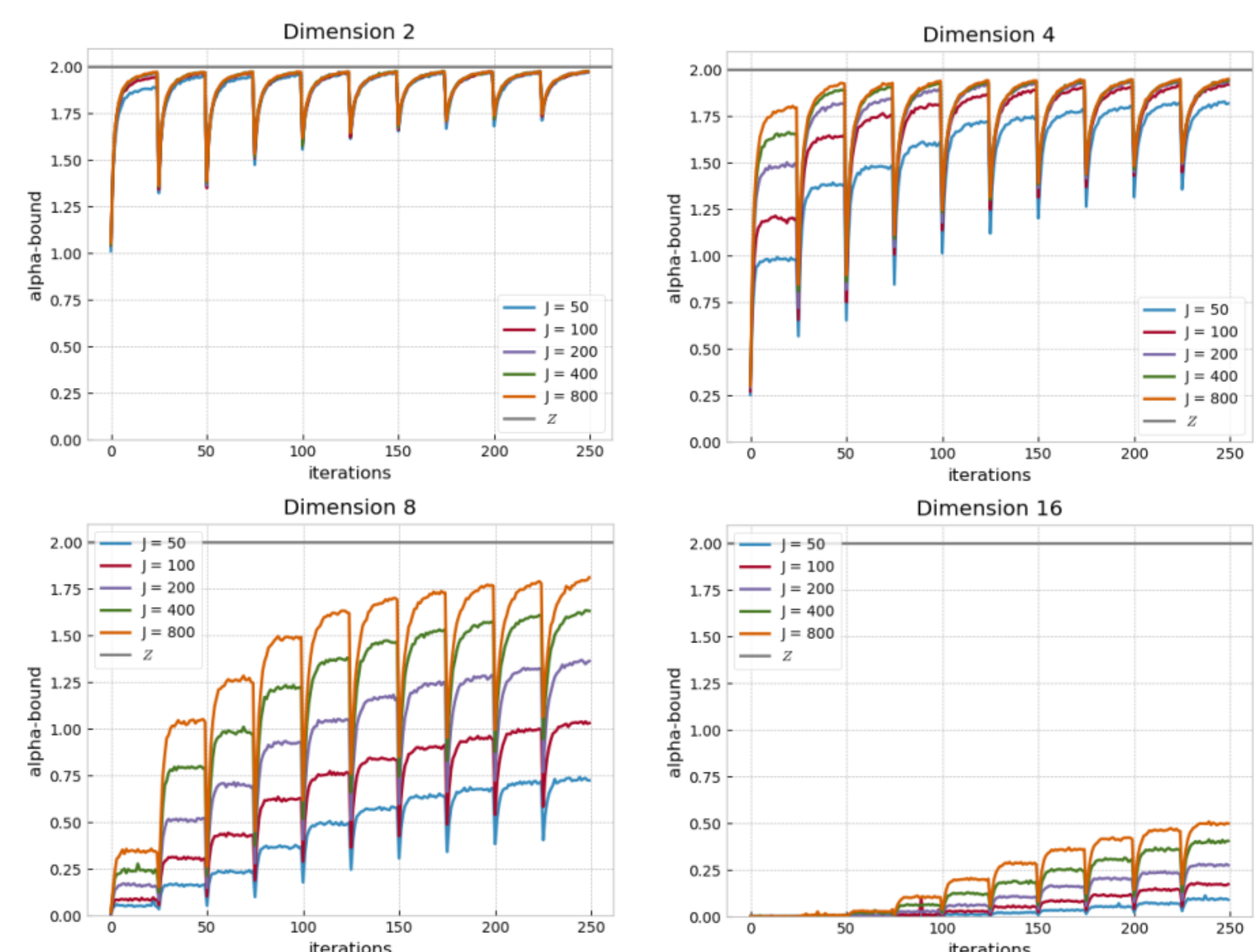
Numerical Experiments

Impact of ϕ (with fixed parameter set $\Theta = \{-2, 2\}$)



$p^*(y) = Z \times [\gamma_1 \mathcal{N}(y; -s, 1) + \gamma_2 \mathcal{N}(y; s, 1)]$, where $\gamma_1 = 0.8$, $\gamma_2 = 0.2$, $s = 2$ and $Z = 2$

Impact of d and J (fully adaptive algorithm)



$p^*(y) = Z \times [0.5 \mathcal{N}(y; -s \mathbf{u}_d, \mathbf{I}_d) + 0.5 \mathcal{N}(y; s \mathbf{u}_d, \mathbf{I}_d)]$ with $s = 2$ and $Z = 2$

References

- [1] R. Douc et al. Convergence of adaptive mixtures of importance sampling schemes. AOS 2007
- [2] JM Hernández-Lobato et al. Black-box α -divergence Minimization. ICML 2015.
- [3] Y. Li and RE. Turner. Rényi Divergence Variational Inference. NIPS 2016.