

Stochastic gradient estimation for discrete latent variables

Rakitin Denis
Higher School of Economics
Bayesian Methods Research Group

Gradient estimation

$$L(\phi) = \mathbb{E}_{q_\phi(z)} f(z)$$

Relaxation-based methods:

- Gumbel-Softmax / Stochastic Softmax Tricks
- REBAR

Score-function methods:

- REINFORCE (+ baselines)
- REBAR/RELAX

AR estimator

Dirichlet reparameterization of $b \sim \text{Cat}(\sigma(\theta))$:

$$\pi \sim \text{Dir}(1_C) \text{ and } b = \arg \min \pi_i e^{-\theta_i}$$

Augmented functional:

$$\mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} f(b) = \mathbb{E}_{\pi \sim \text{Dir}(1_C)} f(\arg \min \pi_i e^{-\theta_i})$$

Augmented true gradient:

$$\nabla_{\theta_l} \mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} f(b) = \mathbb{E}_{\pi} f(b)(1 - C\pi_l) \approx f(b)(1 - C\pi_l)$$

$$\pi \sim \text{Dir}(1_C), \quad b = \arg \min \pi_i e^{-\theta_i}$$

AR estimator

Dirichlet reparameterization of $b \sim \text{Cat}(\sigma(\theta))$:

$$\pi \sim \text{Dir}(1_C) \text{ and } b = \arg \min \pi_i e^{-\theta_i}$$

Augmented functional:

$$\mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} f(b) = \mathbb{E}_{\pi \sim \text{Dir}(1_C)} f(\arg \min \pi_i e^{-\theta_i})$$

Augmented true gradient:

$$\nabla_{\theta_l} \mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} f(b) = \mathbb{E}_{\pi} f(b)(1 - C\pi_l) \approx f(b)(1 - C\pi_l)$$

$$\pi \sim \text{Dir}(1_C), \quad b = \arg \min \pi_i e^{-\theta_i}$$

Score term

AR: what happened

Exponential racing:

$$\tau_1, \dots, \tau_C \sim \text{Exp}(e^{\theta_1}), \dots, \text{Exp}(e^{\theta_C})$$

$$b = \arg \min \tau_i \sim \text{Cat}(\sigma(\theta))$$

Exponential reparameterization:

$$\varepsilon \sim \text{Exp}(1) \text{ then } (1/\lambda) \varepsilon \sim \text{Exp}(\lambda)$$

Dirichlet vs Gamma:

$$(\varepsilon_1, \dots, \varepsilon_C) \stackrel{d}{=} \pi \cdot \xi, \text{ where}$$

$$\pi \sim \text{Dir}(1_C) \text{ and } \xi \sim \Gamma(C, 1)$$

$$\begin{aligned} b &\rightarrow \arg \min \tau_i \rightarrow \arg \min \varepsilon_i e^{-\theta_i} \\ &\rightarrow \arg \min \xi \pi_i e^{-\theta_i} \rightarrow \arg \min \pi_i e^{-\theta_i} \end{aligned}$$

AR: symmetries

Notation:

$$\pi^{j \rightleftharpoons l} = (\pi_1, \dots, \pi_l, \dots, \pi_j, \dots, \pi_C)$$

$$b = \arg \min_i \pi_i e^{-\theta_i}, \quad b^{j \rightleftharpoons l} = \arg \min_i \pi_i^{j \rightleftharpoons l} e^{-\theta_i}$$

Dirichlet symmetry:

$$\pi \sim \text{Dir}(1_C) \text{ then } \pi^{j \rightleftharpoons l} \sim \text{Dir}(1_C)$$

Rewrite the gradient:

$$\nabla_{\theta_l} \mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} f(b) = \mathbb{E}_{\pi} f(b)(1 - C\pi_l) = \mathbb{E}_{\pi} f(b^{j \rightleftharpoons l})(1 - C\pi_j)$$

Reference index j
Differentiate w.r.t. l

AR: deriving a baseline

$$\nabla_{\theta_l} \mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} f(b) = \mathbb{E}_{\pi} f(b)(1 - C\pi_l) = \mathbb{E}_{\pi} f(b^{j \Rightarrow l})(1 - C\pi_j)$$

Let's first write:

$$\begin{aligned} 0 &= \frac{1}{C} \sum_{m=1}^C (1 - C\pi_m) = \frac{1}{C} \sum_{m=1}^C f(b)(1 - C\pi_m) = \\ &= \frac{1}{C} \mathbb{E}_{\pi} \sum_{m=1}^C f(b)(1 - C\pi_m) = \frac{1}{C} \mathbb{E}_{\pi} \sum_{m=1}^C f(b^{j \Rightarrow m})(1 - C\pi_j) \end{aligned}$$

Reference index j
Differentiate w.r.t. l

ARS estimator follows:

$$\mathbb{E}_{\pi} f(b^{j \Rightarrow l})(1 - C\pi_j) = \mathbb{E}_{\pi} \left(f(b^{j \Rightarrow l}) - \frac{1}{C} \sum_{m=1}^C f(b^{j \Rightarrow m}) \right) (1 - C\pi_j)$$

ARS and ARSM

ARS = AR + symmetries + baseline

Number of function evaluations

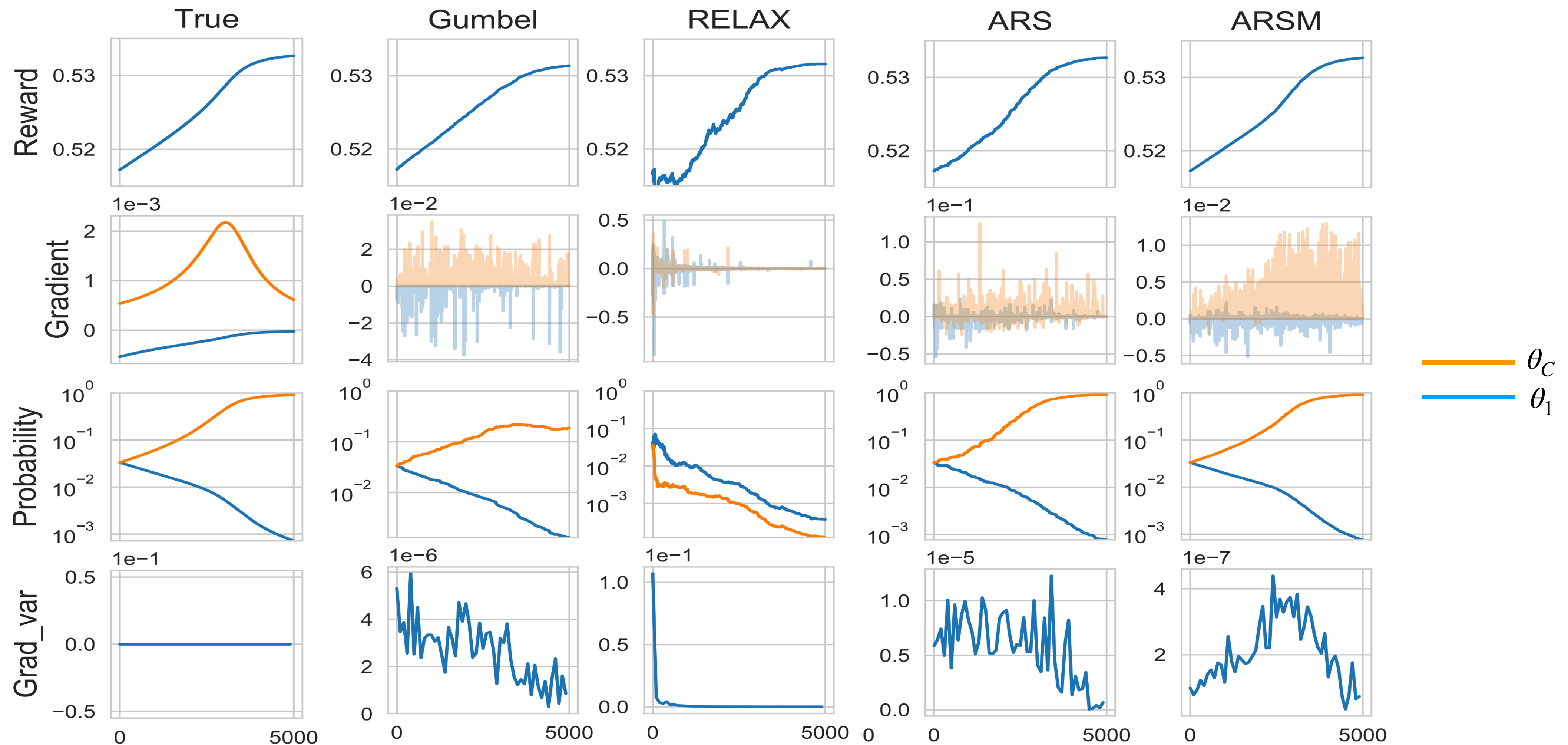
$$\nabla_{\theta_l} L(\theta) = \mathbb{E}_{\pi} \left(\underbrace{f(b^{j \Rightarrow l}) - \frac{1}{C} \sum_{m=1}^C f(b^{j \Rightarrow m})}_{\mathcal{G}_{ARS, l}} \right) (1 - C\pi_j) \leq C$$

ARSM = ARS + mean over reference indices

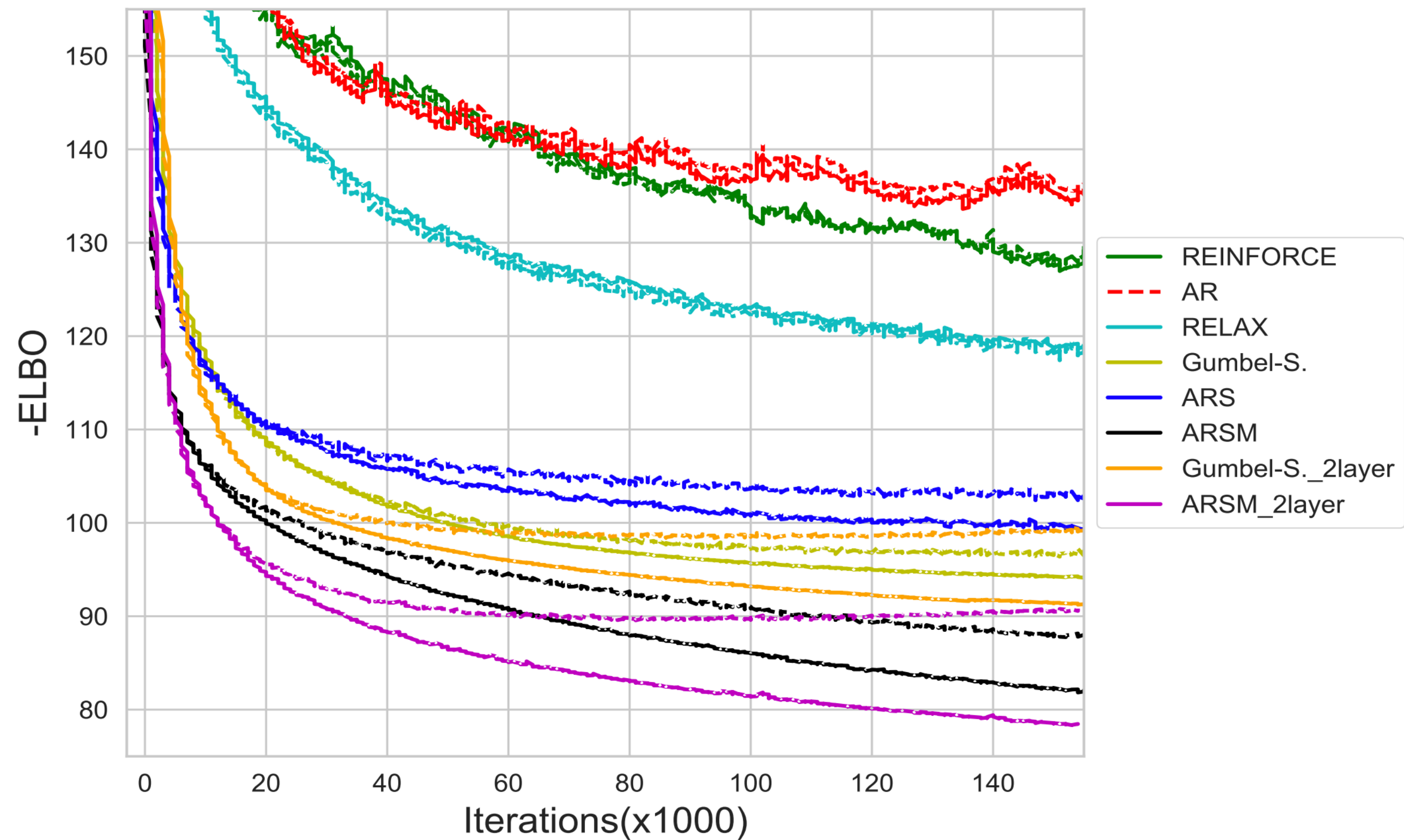
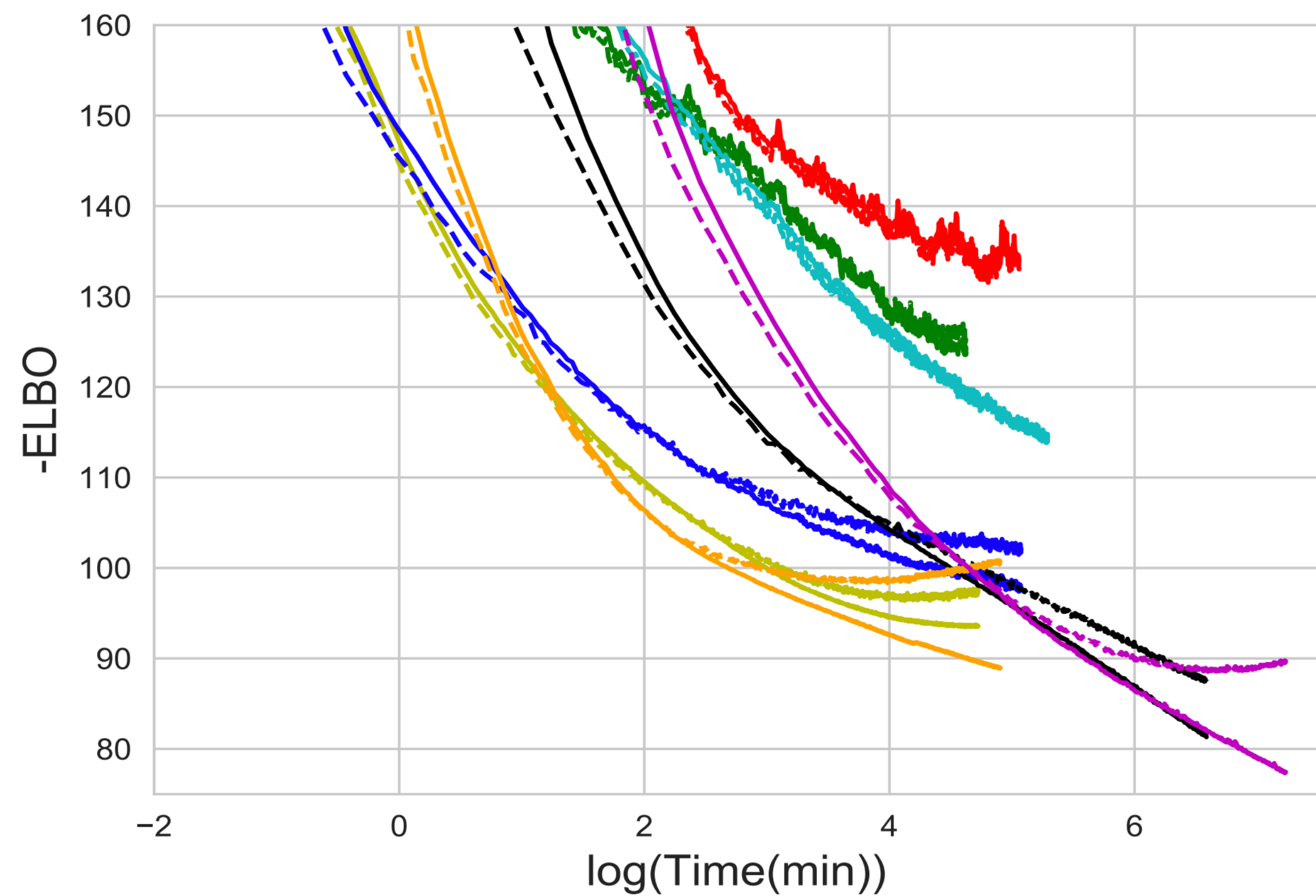
$$\nabla_{\theta_l} L(\theta) = \mathbb{E}_{\pi} \left[\underbrace{\frac{1}{C} \sum_{j=1}^C \left(f(b^{j \Rightarrow l}) - \frac{1}{C} \sum_{m=1}^C f(b^{j \Rightarrow m}) \right) (1 - C\pi_j)}_{\mathcal{G}_{ARSM, l}} \right] \leq C(C-1)/2 + 1$$

Toy experiment

$$L(\theta) = \mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} f(b) = \mathbb{E}_{b \sim \text{Cat}(\sigma(\theta))} 0.5 + b/(CR)$$



VAE on binarized MNIST



ARM+ and ARSM+

Ideal case - Rao-Blackwellization

$$\begin{aligned}\nabla_{\theta_l} L(\theta) &= \mathbb{E}_{\pi} \left(f(b^{j \Rightarrow l}) - \frac{1}{C} \sum_{m=1}^C f(b^{j \Rightarrow m}) \right) (1 - C\pi_j) = \\ &= \mathbb{E}_{b^{j \Rightarrow 1}, \dots, b^{j \Rightarrow C}} \mathbb{E}_{\pi | b^{j \Rightarrow 1}, \dots, b^{j \Rightarrow C}} \left(f(b^{j \Rightarrow l}) - \frac{1}{C} \sum_{m=1}^C f(b^{j \Rightarrow m}) \right) (1 - C\pi_j)\end{aligned}$$

Reference index j
Differentiate w.r.t. l

We only need to compute:

$$\mathbb{E}_{\pi | b^{j \Rightarrow 1}, \dots, b^{j \Rightarrow C}} [\pi_j] = \mathbb{E}_{\pi_{-(out, j)} | b^{j \Rightarrow 1}, \dots, b^{j \Rightarrow C}} \mathbb{E}_{\pi_j | \pi_{-(out, j)}, b^{j \Rightarrow 1}, \dots, b^{j \Rightarrow C}} [\pi_j]$$

Compute analytically

Partial integrating

$$\mathbb{E}_{\pi_j | \pi_{-(out,j)}, b^{j \Rightarrow 1}, \dots, b^{j \Rightarrow C}}[\pi_j]$$

Observe that $b^{j \Rightarrow m}$ imply bounds on π_j :

$$b^{j \Rightarrow m} = k \text{ means } \arg \min_i \pi_i^{j \Rightarrow m} e^{-\theta_i} = k$$

$$\pi_i^{j \Rightarrow m} e^{-\theta_i} \geq \pi_k^{j \Rightarrow m} e^{-\theta_k} \iff \pi_i^{j \Rightarrow m} e^{\theta_k - \theta_i} \geq \pi_k^{j \Rightarrow m}$$

Dirichlet is uniform on the simplex, then conditional of π_j is also uniform

Leveraging symmetry

Multidimensional case:

$$f(b^{j \rightleftharpoons m}) = f(b_1^{j_1 \rightleftharpoons m}, \dots, b_K^{j_1 \rightleftharpoons m})$$

Introduce

$$\delta_k = I\{b_k^{j \rightleftharpoons 1} = \dots = b_k^{j \rightleftharpoons C}\}$$

$$\mathbb{E}_{\pi|\delta_k=1} \left(f(b^{j \rightleftharpoons l}) - \frac{1}{C} \sum_{m=1}^C f(b^{j \rightleftharpoons m}) \right) (1 - C\pi_{k,j}) =$$

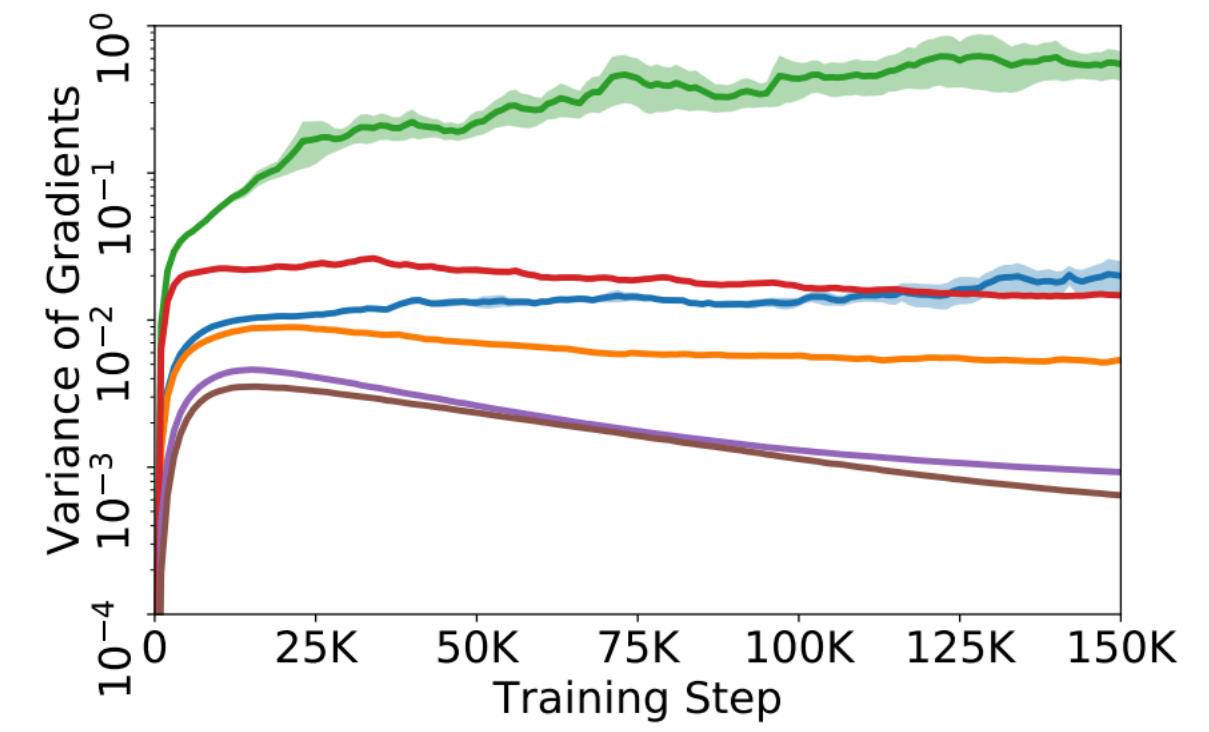
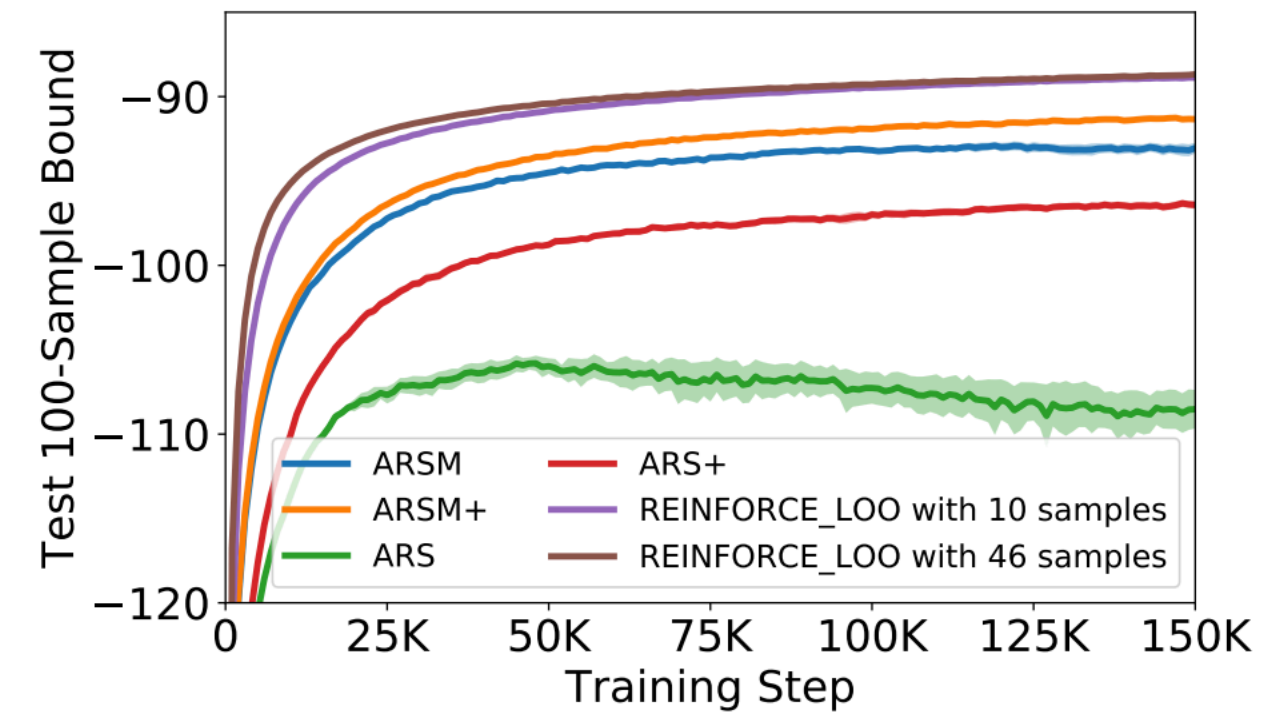
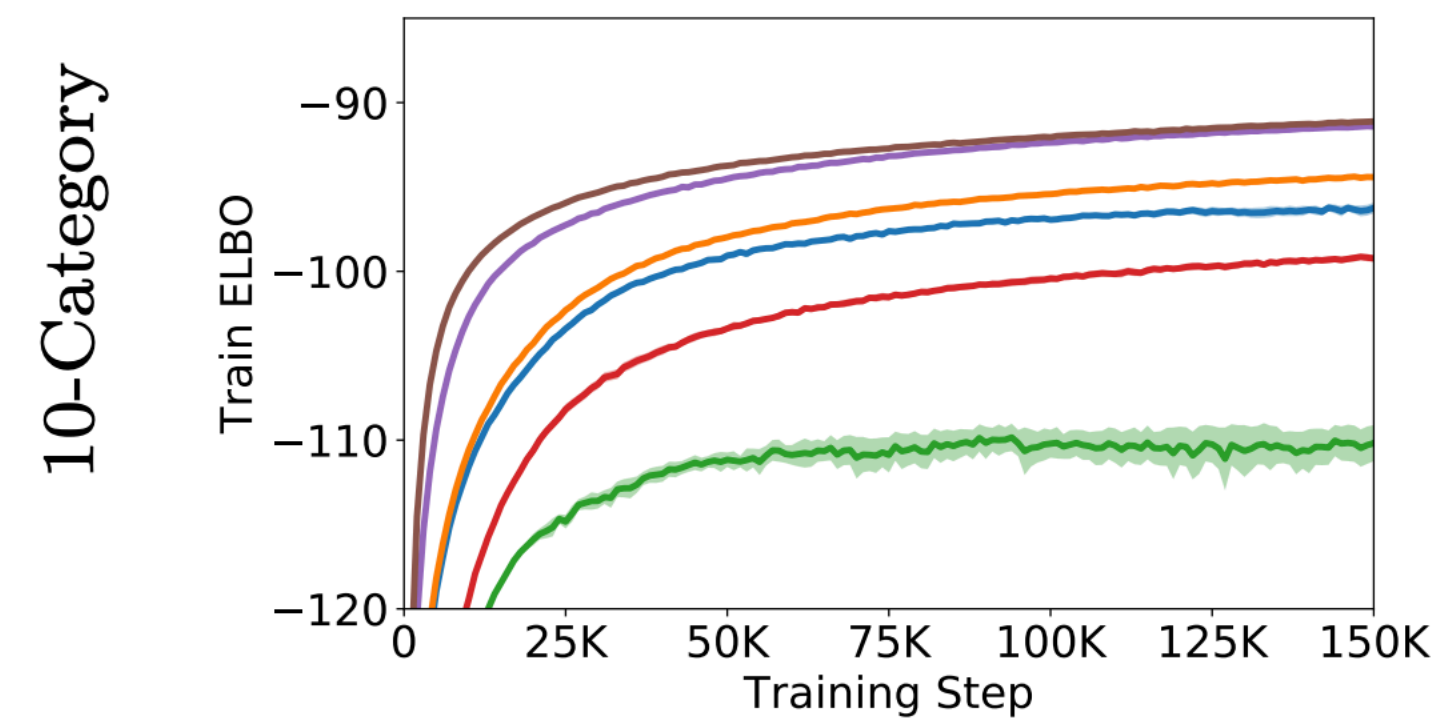
$$= \mathbb{E}_{\pi_k|\delta_k=1} \left(\mathbb{E}_{\pi_{-k}} f(b^{j \rightleftharpoons l}) - \frac{1}{C} \sum_{m=1}^C \mathbb{E}_{\pi_{-k}} f(b^{j \rightleftharpoons m}) \right) = 0$$

Final estimators

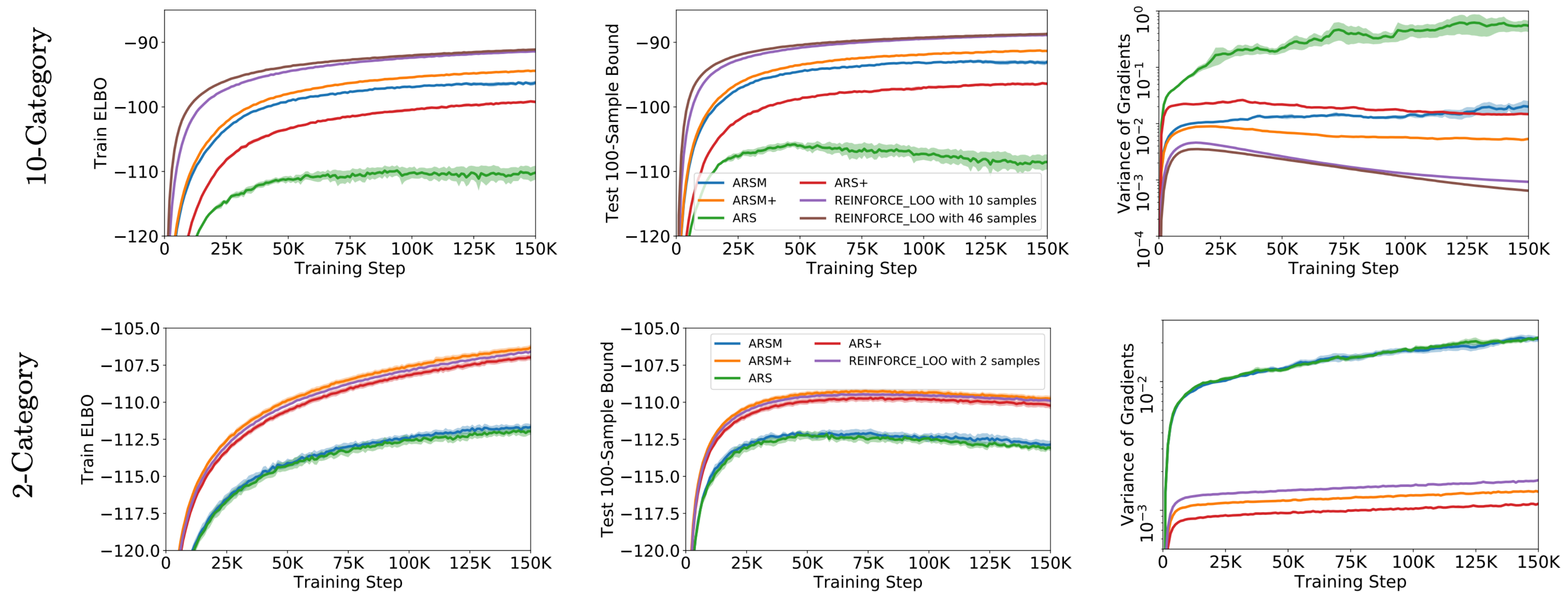
$$g_{ARS+,k,l} = \underbrace{\mathbb{E}_{\pi_{k,j} | \pi_{k,-(1,j)}, b_k^{1 \Rightarrow j}, \dots, b_k^{C \Rightarrow j}}}_{\text{partial integration}} \underbrace{g_{ARS,k,l}(1 - \delta_k)}_{\text{symmetry}}$$

ARSM is modified only by symmetry argument

Experiments



Experiments



REINFORCE+ (LOO)

$$L(\phi) = \mathbb{E}_{q_\phi(z)} f(z)$$

$$g_{RF} = f(z_i) \nabla_\phi \log q_\phi(z_i)$$

$$g_{RF+} = \frac{1}{S} \left(\sum_{i=1}^S \left(f(z_i) - \underbrace{\sum_{j \neq i} f(z_j)}_{\text{loo baseline}} \right) \nabla_\phi \log q_\phi(z_i) \right) = \frac{1}{S-1} \left(\sum_{i=1}^S (f(z_i) - \bar{f}) \nabla_\phi \log q_\phi(z_i) \right)$$

Log-variance loss

$$L(\phi) = ELBO(\phi) = \mathbb{E}_{q_\phi(z)} \log \frac{p(x, z)}{q_\phi(z)}$$

Log-variance loss:

$$\mathcal{L}_r(q_\phi(z) || p(z|x)) = \frac{1}{2} \text{Var}_r \left(\log \frac{q_\phi(z)}{p(z|x)} \right)$$

Gradients property:

$$\nabla_\phi \mathcal{L}_r(q_\phi(z) || p(z|x)) \Big|_{r=q_\phi} = \nabla_\phi KL(q_\phi(z) || p(z|x))$$

Derivation of VarGrad

$$\mathcal{L}_r(q_\phi(z) || p(z|x)) = \frac{1}{2} \text{Var}_r \left(\log \frac{q_\phi(z)}{p(z|x)} \right)$$

Sample variance estimation:

$$\mathcal{L}_r(q_\phi(z) || p(z|x)) \approx \frac{1}{2(S-1)} \sum_{s=1}^S \left(f_\phi(z^{(s)}) - \bar{f}_\phi \right)^2, \quad z^{(s)} \stackrel{\text{i.i.d.}}{\sim} r(z) \quad \left| \text{differentiate + grad property} \right.$$

$$\hat{g}_{\text{VarGrad}}(\phi) = \frac{1}{S-1} \left(\sum_{s=1}^S f_\phi(z^{(s)}) \nabla_\phi \log q_\phi(z^{(s)}) - \bar{f}_\phi \sum_{s=1}^S \nabla_\phi \log q_\phi(z^{(s)}) \right)$$

Unbiased estimate of KL => ELBO

Properties

Another way to introduce RF+:

$$\hat{g}_{\text{CV}}(\phi) = \hat{g}_{\text{Reinforce}}(\phi) - a \odot \left(\frac{1}{S} \sum_{s=1}^S \nabla_{\phi} \log q_{\phi}(z^{(s)}) \right) \quad a = \bar{f}_{\phi} \mathbf{1}$$

Variance minimizing parameter:

$$a_i^* = \frac{\text{Cov}_{q_{\phi}}(f_{\phi} \partial_{\phi_i} \log q_{\phi}, \partial_{\phi_i} \log q_{\phi})}{\text{Var}_{q_{\phi}}(\partial_{\phi_i} \log q_{\phi})}$$

Lemma about mean of parameter in VarGrad:

$$a^* = \mathbb{E}_{q_{\phi}}[a^{\text{VarGrad}}] + \delta^{\text{CV}} = -\text{ELBO}(\phi) + \delta^{\text{CV}} \quad \delta_i^{\text{CV}} = \frac{\text{Cov}_{q_{\phi}}\left(f_{\phi}, (\partial_{\phi_i} \log q_{\phi})^2\right)}{\text{Var}_{q_{\phi}}(\partial_{\phi_i} \log q_{\phi})}$$

Properties

Main lemma: correction term is small

Suppose $\sup_z \frac{q_\phi(z)}{p(z|x)} < C$ and define $\text{Kurt}[\partial_{\phi_i} \log q_\phi] = \frac{\mathbb{E}_{q_\phi}[(\partial_{\phi_i} \log q_\phi)^4]}{(\mathbb{E}_{q_\phi}[(\partial_{\phi_i} \log q_\phi)^2])^2}$

Then

$$\left| \frac{\delta_i^{\text{CV}}}{\mathbb{E}_{q_\phi}[a^{\text{VarGrad}}]} \right| \leq \frac{2\sqrt{C \text{Kurt}[\partial_{\phi_i} \log q_\phi]}}{\left| \sqrt{\text{KL}(q_\phi(z) \parallel p(z|x))} - \frac{\log p(x)}{\sqrt{\text{KL}(q_\phi(z) \parallel p(z|x))}} \right|}$$

KL is large: $\left| \frac{\delta_i^{\text{CV}}}{\mathbb{E}_{q_\phi}[a^{\text{VarGrad}}]} \right| \lesssim \mathcal{O}\left(\text{KL}(q_\phi(z) \parallel p(z|x))^{-1/2}\right)$

KL is small: $\left| \frac{\delta_i^{\text{CV}}}{\mathbb{E}_{q_\phi}[a^{\text{VarGrad}}]} \right| \lesssim \mathcal{O}\left(\text{KL}(q_\phi(z) \parallel p(z|x))^{1/2}\right)$

Properties

VarGrad is better than REINFORCE:

Suppose $-\frac{\delta_i^{CV}}{\mathbb{E}_{q_\phi}[a^{VarGrad}]} = \frac{\delta_i^{CV}}{\text{ELBO}(\phi)} < \frac{1}{2}$, then there exists S_0 such that

$$\text{Var}(\hat{g}_{VarGrad,i}(\phi)) \leq \text{Var}(\hat{g}_{Reinforce,i}(\phi)) \text{ for all } S > S_0$$

Corollary:

If KL grows w.r.t. latent dimension, then:

There exist $S_0, D_0 \in \mathbb{N}$ for which

$$\text{Var}(\hat{g}_{VarGrad,i}(\phi)) \leq \text{Var}(\hat{g}_{Reinforce,i}(\phi)) \text{ for all } S > S_0$$