# ARM: Augment-REINFORCE-merge gradient for stochastic binary networks

Mingzhang Yin, Mingyuan Zhou

Department of Statistics and Data Sciences
The University of Texas at Austin

October 12, 2018

# Discrete latent variables

- Discrete latent variables are widely used in mixture models, mixed membership model, variable selection, and sigmoid belief net etc.

- A common task is to optimize

$$\mathcal{E}(\phi) = \int f(\mathbf{z}) q_\phi(\mathbf{z}) d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})}[f(\mathbf{z})]$$

- This objective includes
  - Maximizing the marginal likelihood of a hierarchical Bayesian model
  - Maximizing the evidence lower bound (ELBO) in variational inference
  - Maximizing the expected reward in reinforcement learning
  - ...

# Unbiased gradient estimators

**Reparameterization**

If $\nabla_z f(z)$ is tractable to compute and $z \sim q_\phi(z)$ can be generated via reparameterization as $z = \mathcal{T}_\phi(\epsilon)$, $\epsilon \sim p(\epsilon)$, then one may apply the reparameterization trick

$$\nabla_\phi \mathcal{E}(\phi) = \nabla_\phi \mathbb{E}_{\epsilon \sim p(\epsilon)}[f(\mathcal{T}_\phi(\epsilon))] = \mathbb{E}_{\epsilon \sim p(\epsilon)}[\nabla_\phi f(\mathcal{T}_\phi(\epsilon))]$$

**REINFORCE (score-function estimator)**

If $\mathbb{E}_{z \sim q_\phi(z)}[\nabla_\phi f(z)] = 0$, using the score function
$\nabla_\phi \log q_\phi(z) = \nabla_\phi q_\phi(z) / q_\phi(z)$, one may use REINFORCE as

$$\nabla_\phi \mathcal{E}(\phi) = \mathbb{E}_{z \sim q_\phi(z)}[f(z) \nabla_\phi \log q_\phi(z)] \approx \frac{1}{K} \sum_{k=1}^{K} f(z^{(k)}) \nabla_\phi \log q_\phi(z^{(k)})$$

## Problems

However, neither estimator is problem free:

- The reparameterization trick requires $f(z)$ to be differentiable and cannot be applied to discrete $z$
- REINFORCE suffers from high Monte Carlo estimation variance

# Gradient estimation for discrete latent variables

For discrete latent variable $z$, to compute the gradient of $\mathcal{E}(\phi) = \mathbb{E}_{z \sim q_\phi(z)}[f(z)]$, existing solutions include

- Biased but low-variance gradient estimator via a continuous relaxation of discrete random variables
  - Gumbel-softmax trick, $\operatorname{argmax}(g_i + \log \pi_i)$

- Variance reduction by adding control variates (a.k.a. baselines)

$$\nabla_\phi \mathcal{E}(\phi) = \mathbb{E}_{z \sim q_\phi(z)}[f(z)\nabla_\phi \log q_\phi(z) - b(z)]$$

where $\mathbb{E}_{z \sim q_\phi(z)}[b(z)] = 0$, for example
  - $b(z) = m\nabla_\phi \log q_\phi(z)$
  - Learnable control variate function parameterized with neural network (REBAR/RELAX)

# ARM gradient

---

**Algorithm 1:** ARM gradient for $V$-dimensional binary latent vector

**input** : Bernoulli distribution $\{q_{\phi_v}(z_v)\}_{v=1:V}$ with probability $\{\sigma(\phi_v)\}_{v=1:V}$, target $f(z)$;
$\quad\quad z = (z_1, \cdots, z_V)$, $\phi = (\phi_1, \cdots, \phi_V)$
**output :** $\phi$ and $\psi$ that maximize $\mathcal{E}(\phi, \psi) = \mathbb{E}_{z \sim \prod_{v=1}^V q_{\phi_v}(z_v)}[f(z; \psi)]$

Initialize $\phi$, $\psi$ randomly;
**while** *not converged* **do**
$\quad$ Sample $z_v \sim$ Bernoulli$(\sigma(\phi_v))$ for $v = 1, \cdots, V$ ;
$\quad$ sample $u_v \sim$ Uniform$(0, 1)$ for $v = 1, \cdots, V$, $u = (u_1, \cdots, u_V)$ ;
$\quad$ $g_\psi = \nabla_\psi f(z; \psi)$ ;
$\quad$ $f_\Delta(u, \phi) = f\big(\mathbf{1}_{[u > \sigma(-\phi)]}\big) - f\big(\mathbf{1}_{[u < \sigma(\phi)]}\big)$ ;
$\quad$ $g_\phi = f_\Delta(u, \phi)(u - \frac{1}{2})$
$\quad$ $\phi = \phi + \rho_t g_\phi$, $\quad \psi = \psi + \eta_t g_\psi$ $\quad$ with step-size $\rho_t$, $\eta_t$
**end**

---

## Augmentation of binary random variable

Bernoulli random variable $z \sim \text{Bernoulli}(\sigma(\phi))$ can be reparameterized by racing two augmented exponential random variables as

$$z = \mathbf{1}_{[\epsilon_1 < \epsilon_2 e^\phi]}, \ \epsilon_1 \sim \text{Exp}(1), \ \epsilon_2 \sim \text{Exp}(1).$$

The objective can be rewritten with respect to exponential random variables as

$$\mathcal{E}(\phi) = \mathbb{E}_{z \sim \text{Bernoulli}(\sigma(\phi))}[f(z)] = \mathbb{E}_{\epsilon_1, \epsilon_2 \overset{iid}{\sim} \text{Exp}(1)}[f(\mathbf{1}_{[\epsilon_1 e^{-\phi} < \epsilon_2]})].$$

## REINFORCE in augmented space

Use REINFORCE with $t_1 = \epsilon_1 e^{-\phi} \sim \mathsf{Exp}(e^{\phi})$ we have

$$
\begin{aligned}
\nabla_\phi \mathcal{E}(\phi) &= \mathbb{E}_{t_1 \sim \mathsf{Exp}(e^\phi), \ \epsilon_2 \sim \mathsf{Exp}(1)}[f(\mathbf{1}_{[t_1 < \epsilon_2]}) \nabla_\phi \log \mathsf{Exp}(t_1; e^\phi)] \\
&= \mathbb{E}_{t_1 \sim \mathsf{Exp}(e^\phi), \ \epsilon_2 \sim \mathsf{Exp}(1)}[f(\mathbf{1}_{[t_1 < \epsilon_2]})(1 - t_1 e^\phi)] \\
&= \mathbb{E}_{\epsilon_1, \epsilon_2 \overset{iid}{\sim} \mathsf{Exp}(1)}[f(\mathbf{1}_{[\epsilon_1 e^{-\phi} < \epsilon_2]})(1 - \epsilon_1)]
\end{aligned}
$$

Similarly with $t_2 = \epsilon_2 e^\phi \sim \mathsf{Exp}(e^{-\phi})$

$$
\begin{aligned}
\nabla_\phi \mathcal{E}(\phi) &= \mathbb{E}_{\epsilon_1 \sim \mathsf{Exp}(1), \ t_2 \sim \mathsf{Exp}(e^{-\phi})}[f(\mathbf{1}_{[\epsilon_1 < t_2]}) \nabla_\phi \log \mathsf{Exp}(t_2; e^{-\phi})] \\
&= -\mathbb{E}_{\epsilon_1 \sim \mathsf{Exp}(1), \ t_2 \sim \mathsf{Exp}(e^{-\phi})}[f(\mathbf{1}_{[\epsilon_1 < t_2]})(1 - t_2 e^{-\phi})] \\
&= -\mathbb{E}_{\epsilon_1, \epsilon_2 \overset{iid}{\sim} \mathsf{Exp}(1)}[f(\mathbf{1}_{[\epsilon_1 e^{-\phi} < \epsilon_2]})(1 - \epsilon_2)]
\end{aligned}
$$

## Merge the gradients

Swapping the indices of the two *iid* standard exponential random variables, the gradient $\nabla_\phi \mathcal{E}(\phi)$ can be equivalently expressed as

$$\nabla_\phi \mathcal{E}(\phi) = -\mathbb{E}_{\epsilon_1, \epsilon_2 \overset{iid}{\sim} \mathsf{Exp}(1)} [f(\mathbf{1}_{[\epsilon_2 e^{-\phi} < \epsilon_1]})(1 - \epsilon_1)].$$

Taking the average leads to

$$\nabla_\phi \mathcal{E}(\phi) = \mathbb{E}_{\epsilon_1, \epsilon_2 \overset{iid}{\sim} \mathsf{Exp}(1)} \left[ \left( f(\mathbf{1}_{[\epsilon_1 e^{-\phi} < \epsilon_2]}) - f(\mathbf{1}_{[\epsilon_2 e^{-\phi} < \epsilon_1]}) \right)(1/2 - \epsilon_1/2) \right]$$

## Merge the gradients

Letting $\epsilon_1, \epsilon_2 \overset{iid}{\sim} \text{Exp}(1)$ is the same in distribution as letting

$$\epsilon_1 = \epsilon u, \ \epsilon_2 = \epsilon(1-u), \quad \text{where } u \sim \text{Uniform}(0,1), \ \epsilon \sim \text{Gamma}(2,1)$$

Marginalizing Gamma R.V. out reduces variance by the law of total variance

$$\begin{aligned}
\nabla_\phi \mathcal{E}(\phi) &= \mathbb{E}_{u \sim \text{U}(0,1), \ \epsilon \sim \text{Ga}(2,1)} \left[ \left( f(\mathbf{1}_{[u>\sigma(-\phi)]}) - f(\mathbf{1}_{[u<\sigma(\phi)]}) \right) (\epsilon u/2 - 1/2) \right], \\
&= \mathbb{E}_{u \sim \text{Uniform}(0,1)} \left[ \left( f(\mathbf{1}_{[u>\sigma(-\phi)]}) - f(\mathbf{1}_{[u<\sigma(\phi)]}) \right) (u - 1/2) \right]
\end{aligned}$$

## Multivariate generalization

Apply the univariate result to the inner expectation

$$\nabla_{\phi_v}\mathcal{E}(\phi) = \mathbb{E}_{\mathbf{z}_{\backslash v}\sim\prod_{\nu\neq v}\text{Bernoulli}(z_\nu;\sigma(\phi_\nu))}\{\nabla_{\phi_v}\mathbb{E}_{z_v\sim\text{Bernoulli}(\sigma(\phi_v))}[f(\mathbf{z})]\}$$
$$= \mathbb{E}_{\mathbf{z}_{\backslash v}\sim\prod_{\nu\neq v}\text{Bernoulli}(z_\nu;\sigma(\phi_\nu))}\{\mathbb{E}_{u_v\sim\text{Uniform}(0,1)}\big[(u_v - 1/2)$$
$$\times\big(f(\mathbf{z}_{\backslash v}, z_v = \mathbf{1}_{[u_v>\sigma(-\phi_v)]}) - f(\mathbf{z}_{\backslash v}, z_v = \mathbf{1}_{[u_v<\sigma(\phi_v)]}))\big)\big]\}.$$

then sample $\mathbf{z}_{\backslash v}$ as $\mathbf{1}_{[u_{\backslash v}<\sigma(\phi_{\backslash v})]}$ or as $\mathbf{z}_{\backslash v} = \mathbf{1}_{[u_{\backslash v}>\sigma(-\phi_{\backslash v})]}$,

$$\nabla_{\phi_v}\mathcal{E}(\phi) = \mathbb{E}_{u_v\sim\text{Uniform}(0,1)}\big\{(u_v - 1/2)\,\mathbb{E}_{\mathbf{z}_{\backslash v}\sim\prod_{\nu\neq v}\text{Bernoulli}(z_\nu;\sigma(\phi_\nu))}\big[$$
$$f(\mathbf{z}_{\backslash v}, z_v = \mathbf{1}_{[u_v>\sigma(-\phi_v)]}) - f(\mathbf{z}_{\backslash v}, z_v = \mathbf{1}_{[u_v<\sigma(\phi_v)]})\big]\big\}$$
$$= \mathbb{E}_{\mathbf{u}\sim\prod_{v=1}^V\text{Uniform}(u_v;0,1)}\big[(u_v - 1/2)\big(f(\mathbf{1}_{[\mathbf{u}>\sigma(-\phi)]}) - f(\mathbf{1}_{[\mathbf{u}<\sigma(\phi)]}))\big]$$

## Antithetic Sampling

The merging step is equivalent to the antithetic sampling

$$\nabla_{\phi_v} \mathcal{E}(\phi) = \mathbb{E}_{\boldsymbol{z}_{\backslash v} \sim \prod_{\nu \neq v} \mathsf{Bernoulli}(z_\nu; \sigma(\phi_\nu))} \{ \nabla_{\phi_v} \mathbb{E}_{z_v \sim \mathsf{Bernoulli}(\sigma(\phi_v))}[f(\boldsymbol{z})] \}$$
$$= \mathbb{E}_{\boldsymbol{u} \sim \prod_{v=1}^{V} \mathsf{Uniform}(u_v; 0,1)} \big[ (1 - 2u_v) f(\mathbf{1}_{[\boldsymbol{u} < \sigma(\phi)]}) \big],$$

With the antithetic sample $\tilde{\boldsymbol{u}} = 1 - \boldsymbol{u}$

$$\nabla_\phi \mathcal{E}(\phi) = \mathbb{E}_{\boldsymbol{u} \sim \prod_v \mathsf{Uniform}(0,1)}[f(\mathbf{1}_{[\boldsymbol{u} < \sigma(\phi)]})(1/2 - \boldsymbol{u})]$$
$$+ \mathbb{E}_{\tilde{\boldsymbol{u}} \sim \prod_v \mathsf{Uniform}(0,1)}[f(\mathbf{1}_{[\tilde{\boldsymbol{u}} < \sigma(\phi)]})(1/2 - \tilde{\boldsymbol{u}})]$$
$$= \mathbb{E}_{\boldsymbol{u} \sim \prod_v \mathsf{Uniform}(0,1)}[f(\mathbf{1}_{[\boldsymbol{u} < \sigma(\phi)]})(1/2 - \boldsymbol{u}) + f(\mathbf{1}_{[\tilde{\boldsymbol{u}} < \sigma(\phi)]})(1/2 - \tilde{\boldsymbol{u}})]$$

# Relation to control variates

### Proposition

Let $\mathcal{B} = \{b(\boldsymbol{u}) : b(\boldsymbol{u}) = -b(1 - \boldsymbol{u}), \}$, and $g_{AR}(\boldsymbol{u}) = f(\mathbf{1}_{[\boldsymbol{u} < \sigma(\phi)]})(1 - 2\boldsymbol{u})$, then with

$$b^*(\boldsymbol{u}) = \arg\min_{b(\boldsymbol{u}) \in \mathcal{B}} Var[g_{AR,v}(\boldsymbol{u}) - b_v(\boldsymbol{u})]$$

we have $g_{ARM}(\boldsymbol{u}) = g_{AR}(\boldsymbol{u}) - b^*(\boldsymbol{u})$ with guaranteed variance reduction.

The common score baseline $b(\boldsymbol{u}) = m(1 - 2\boldsymbol{u}) \in \mathcal{B}$, which means the variance reduction of ARM gradient is better than score baseline even with optimal $m$.

# A simple example

Learning $\phi$ to maximize $\mathcal{E}(\phi) = \mathbb{E}_{z \sim \text{Bernoulli}(\sigma(\phi))}[(z - p_0)^2]$, where $p_0 \in \{0.49, 0.499, 0.501, 0.51\}$
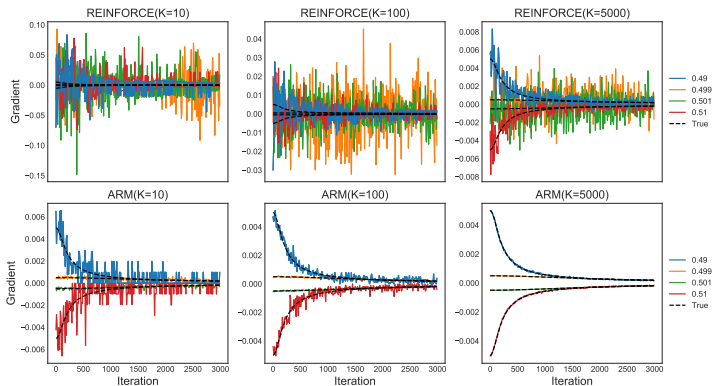
# A simple example



Figure: Estimation of the true gradient at each iteration using $K > 1$ Monte Carlo samples, using REINFORCE, shown in the top row, or ARM, shown in the bottom row. The ARM estimator exhibits significant lower variance given the same number of Monte Carlo samples.

# Multiple discrete stochastic layers

- A latent variable model with multiple stochastic hidden layers is

$$\boldsymbol{x} \sim p_{\theta_0}(\boldsymbol{x} \,|\, \boldsymbol{b}_1), \ \boldsymbol{b}_1 \sim p_{\theta_1}(\boldsymbol{b}_1 \,|\, \boldsymbol{b}_2), \ldots, \boldsymbol{b}_T \sim p_{\theta_T}(\boldsymbol{b}_T),$$

- The joint likelihood is

$$p(\boldsymbol{x}, \boldsymbol{b}_{1:T} \,|\, \boldsymbol{\theta}_{0:T}) = p_{\theta_0}(\boldsymbol{x} \,|\, \boldsymbol{b}_1) \Big[ \prod_{t=1}^{T-1} p_{\theta_t}(\boldsymbol{b}_t \,|\, \boldsymbol{b}_{t+1}) \Big] p_{\theta_T}(\boldsymbol{b}_T).$$

# VAE with multiple discrete stochastic layers

- VAE optimizes ELBO as $\mathbb{E}_{q_{\boldsymbol{w}}(\boldsymbol{b}|\boldsymbol{x})} \log \frac{p_\theta(\boldsymbol{x}|\boldsymbol{b})p(\boldsymbol{b})}{q_{\boldsymbol{w}}(\boldsymbol{b}|\boldsymbol{x})}$
- The encoder is designed as

$$q_{\boldsymbol{w}_{1:T}}(\boldsymbol{b}_{1:T} \mid \boldsymbol{x}) = q_{\boldsymbol{w}_1}(\boldsymbol{b}_1 \mid \boldsymbol{x}) \Big[ \prod_{t=1}^{T-1} q_{\boldsymbol{w}_{t+1}}(\boldsymbol{b}_{t+1} \mid \boldsymbol{b}_t) \Big]$$

- $q_{\boldsymbol{w}_t}(\boldsymbol{b}_t \mid \boldsymbol{b}_{t-1}) = \text{Bernoulli}(\boldsymbol{b}_t; \sigma(\mathcal{T}_{\boldsymbol{w}_t}(\boldsymbol{b}_{t-1})))$
- The ELBO can be expressed as

$$\mathcal{E}(\boldsymbol{w}_{1:T}) = \mathbb{E}_{\boldsymbol{b}_{1:T} \sim q_{\boldsymbol{w}_{1:T}}(\boldsymbol{b}_{1:T} \mid \boldsymbol{x})} \left[ f(\boldsymbol{b}_{1:T}) \right], \text{ where}$$
$$f(\boldsymbol{b}_{1:T}) = \log p_{\theta_0}(\boldsymbol{x} \mid \boldsymbol{b}_1) + \log p_{\theta_{1:T}}(\boldsymbol{b}_{1:T}) - \log q_{\boldsymbol{w}_{1:T}}(\boldsymbol{b}_{1:T} \mid \boldsymbol{x}).$$

## ARM gradient of the ELBO

To compute the gradient with respect to $\boldsymbol{w}_t$,

$$\mathcal{E}(\boldsymbol{w}_{1:T}) = \mathbb{E}_{q(\boldsymbol{b}_{1:t-1})}\mathbb{E}_{q(\boldsymbol{b}_t \mid \boldsymbol{b}_{t-1})}\mathbb{E}_{q(\boldsymbol{b}_{t+1:T} \mid \boldsymbol{b}_t)}[f(\boldsymbol{b}_{1:T})]$$

we treat $\mathbb{E}_{q(\boldsymbol{b}_{t+1:T} \mid \boldsymbol{b}_t)}[f(\boldsymbol{b}_{1:T})]$ as a function

$$\nabla_{\boldsymbol{w}_t}\mathcal{E}(\boldsymbol{w}_{1:T})$$
$$= \mathbb{E}_{q(\boldsymbol{b}_{1:t-1})}\left[\mathbb{E}_{\boldsymbol{u}_t \sim \text{Uniform}(0,1)}[f_\Delta(\boldsymbol{u}_t, \mathcal{T}_{\boldsymbol{w}_t}(\boldsymbol{b}_{t-1}), \boldsymbol{b}_{1:t-1})(\boldsymbol{u}_t - 1/2)]\nabla_{\boldsymbol{w}_t}\mathcal{T}_{\boldsymbol{w}_t}(\boldsymbol{b}_{t-1})\right],$$

where

$$f_\Delta(\boldsymbol{u}_t, \mathcal{T}_{\boldsymbol{w}_t}(\boldsymbol{b}_{t-1}), \boldsymbol{b}_{1:t-1}) = \mathbb{E}_{\boldsymbol{b}_{t+1:T} \sim q(\boldsymbol{b}_{t+1:T} \mid \boldsymbol{b}_t),\ \boldsymbol{b}_t = \mathbf{1}[\boldsymbol{u}_t > \sigma(-\mathcal{T}_{\boldsymbol{w}_t}(\boldsymbol{b}_{t-1}))])}[f(\boldsymbol{b}_{1:T})]$$
$$- \mathbb{E}_{\boldsymbol{b}_{t+1:T} \sim q(\boldsymbol{b}_{t+1:T} \mid \boldsymbol{b}_t),\ \boldsymbol{b}_t = \mathbf{1}[\boldsymbol{u}_t < \sigma(\mathcal{T}_{\boldsymbol{w}_t}(\boldsymbol{b}_{t-1}))])}[f(\boldsymbol{b}_{1:T})]$$

# ARM for VAE

| Linear | | Nonlinear | | Two layers | |
|---|---|---|---|---|---|
| Algrithm | $-\log p(x)$ | Algrithm | $-\log p(x)$ | Algrithm | $-\log p(x)$ |
| NVIL [*] | $= 113.1$ | NVIL [*] | $= 102.2$ | NVIL [*] | $= 99.8$ |
| WS [*] | $= 120.8$ | WS [*] | - | WS [*] | $= 107.7$ |
| REBAR [⋆] | $= 107.7$ | REBAR [⋆] | $= 101.4$ | REBAR [⋆] | $= \mathbf{95.4}$ |
| MuProp [†] | $\leq 113.0$ | MuProp [⋆] | $= 99.1$ | MuProp [†] | $\leq 100.4$ |
| Concrete [⋆] | $= 107.3$ | Concrete [⋆] | $= 99.6$ | Concrete [⋆] | $= 95.6$ |
| RELAX [‡] | $\leq 113.6$ | RELAX [‡] | $\leq 119.2$ | RELAX [‡] | $\leq 100.9$ |
| LeGrad | $\leq 117.5$ | LeGrad | - | LeGrad | - |
| ARM | $= \mathbf{107.2 \pm 0.1}$ | ARM | $= \mathbf{98.4 \pm 0.3}$ | ARM | $= 96.7 \pm 0.3$ |

Table 2: Negative log-likelihood of discrete VAEs on test set trained with mutiple stochastic gradient estimators on MNIST where ∗, ⋆, †, ‡ represent the reported results from Mnih & Gregor (2014), Tucker et al. (2017), Gu et al. (2016), Grathwohl et al. (2018) respectively. LeGrad result is obtained by running the author provided code. ARM rerult is over 5 random initializations with standard error of the mean.

| Linear | | Nonlinear | | Two layers | |
|---|---|---|---|---|---|
| Algrithm | $-\log p(x)$ | Algrithm | $-\log p(x)$ | Algrithm | $-\log p(x)$ |
| NVIL [*] | $= 117.6$ | NVIL [*] | $= \mathbf{116.6}$ | NVIL [*] | $= 111.4$ |
| REBAR [⋆] | $= 117.7$ | REBAR [⋆] | $= 118.0$ | REBAR [⋆] | $= 110.8$ |
| MuProp [⋆] | $= 117.6$ | MuProp [⋆] | $= 117.5$ | MuProp [⋆] | $= 111.2$ |
| Concrete [⋆] | $= 117.7$ | Concrete [⋆] | $= 116.7$ | Concrete [⋆] | $= 111.3$ |
| RELAX [‡] | $\leq 122.1$ | RELAX [‡] | $\leq 128.2$ | RELAX [‡] | $\leq 115.4$ |
| ARM | $= \mathbf{115.8 \pm 0.2}$ | ARM | $= 117.6 \pm 0.4$ | ARM | $= \mathbf{109.8 \pm 0.3}$ |

Table 3: Negative log-likelihood of discrete VAEs on test set trained with mutiple stochastic gradient estimators on OMNIGLOT where ∗, ⋆, ‡ represent the reported results from Mnih & Gregor (2014), Tucker et al. (2017), Grathwohl et al. (2018) respectively.
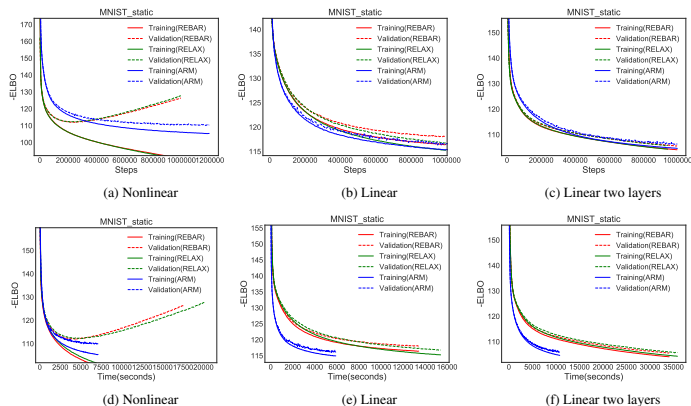
# ARM for VAE



Figure 5: Test negative ELBOs on MNIST(static) with respect to training iterations, shown in the top row, and wall clock times on Tesla-K40 GPU, shown in the bottom row, for three differently structured Bernoulli VAEs.
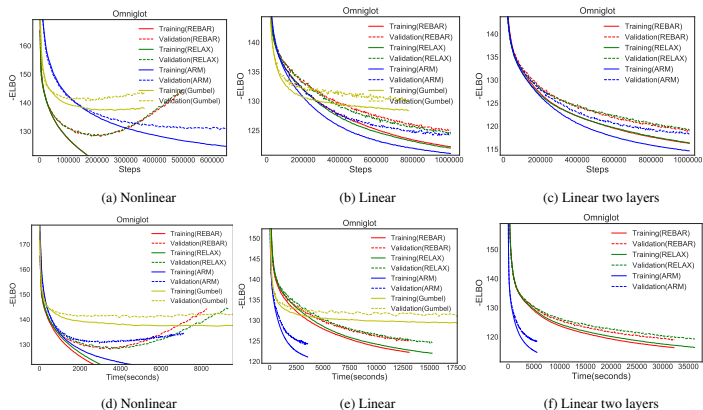
# ARM for VAE



Figure 6: Test negative ELBOs on OMNIGLOT with respect to training iterations, shown in the top row, and wall clock times on Tesla-K40 GPU, shown in the bottom row, for three differently structured Bernoulli VAEs.

## MLE with multiple discrete stochastic layers

- Predict lower half of MNIST digit $x_l$ given the upper half $x_u$;
- Maximizing the conditional likelihood $p_{\theta_{0:2}}(x_l \mid x_u)$
- Approximate $\log p_{\theta_{0:2}}(x_l \mid x_u)$ with

$$\log \frac{1}{K} \sum_{k=1}^{K} \text{Bernoulli}(x_l; \sigma(\mathcal{T}_{\theta_0}(b_1^{(k)})))$$

where $b_1^{(k)} \sim \text{Bernoulli}(\sigma(\mathcal{T}_{\theta_1}(b_2^{(k)})))$, $b_2^{(k)} \sim \text{Bernoulli}(\sigma(\mathcal{T}_{\theta_2}(x_u)))$.

- Training with K=1 and approximate negative log-likelihood $-\log p_{\theta_{0:2}}(x_l \mid x_u)$ with $K = 1000$
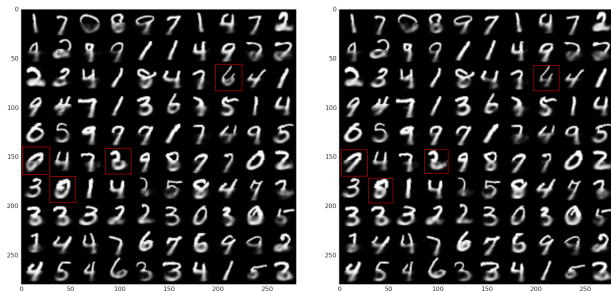
# MLE with multiple discrete stochastic layers



Figure 4: Randomly selected example results of predicting the lower half of a MNIST digit given its upper half, using a binary stochastic network, which has two binary linear stochastic hidden layers and is trained by ARM maximum likelihood inference.

Table 4: For the MNIST conditional distribution estimation benchmark task, comparison of the test negative log-likelihood between various gradient estimators in Jang et al. (2017) is reported here.

| Gradient estimator | ARM | ST | DARN | Annealed ST | ST Gumbel-S. | SF | MuProp |
|---|---|---|---|---|---|---|---|
| $-\log p(\boldsymbol{x}_l \mid \boldsymbol{x}_u)$ | $\mathbf{57.9 \pm 0.1}$ | 58.9 | 59.7 | 58.7 | 59.3 | 72.0 | 58.9 |

# Future work

- ARM gradient can be generalized to categorical variables by the the exponential racing property: $x_i \sim \text{Exp}(\lambda_i)$ are independent exponential random variables for $i = 1, \ldots, M$, then

$$P\big(i = \arg\min_j x_j\big) = P\left(x_i < x_j, \ \forall \ j \neq i\right) = \lambda_i \bigg/ \sum_{i=1}^{M} \lambda_i \ .$$

- Combination with other variance reduction methods
- Applications in the reinforcement learning with discrete actions and natural language processing

# Reference

- Connectionist learning of belief networks Neal (1992)
- Neural Variational Inference and Learning in Belief Networks Mnih and Gregor (2014)
- Local expectation gradients for black box variational inference Titsias and Lázaro-Gredilla (2015)
- MuProp: Unbiased backpropagation for stochastic neural networks Gu et al. (2016)
- Categorical reparameterization with Gumbel-softmaxMaddison et al. (2017); Jang et al. (2017)
- REBAR: Low-variance, unbiased gradient estimates for discrete latent variable models Tucker et al. (2017)
- ARM: Augment-REINFORCE-merge gradient for discrete latent variable models Yin and Zhou (2018)

Thank you!