ARM: Augment-REINFORCE-merge gradient for discrete latent variable models

Mingyuan Zhou* § Joint work with Mingzhang Yin§

*IROM Department, McCombs School of Business [§]Department of Statistics and Data Sciences The University of Texas at Austin

Institut de Recherche en Informatique de Toulouse Toulouse, July 5, 2018

- 4 同 1 - 4 三 1 - 4 三 1

Joint work with



Mingzhang Yin

PhD student (since Fall 2015) in Statistics and Data Sciences

4 E

A B A B A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A

- Discrete latent variables are widely used in mixture models, mixed membership model, sparse factor model, variable selection, etc.
- A common task is to optimize

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

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

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

$$abla_{\phi} \mathcal{E}(\phi) = \mathbb{E}_{oldsymbol{z} \sim q_{\phi}(oldsymbol{z})}[f(oldsymbol{z})
abla_{\phi} \log q_{\phi}(oldsymbol{z})] pprox rac{1}{K} \sum_{k=1}^{K} f(oldsymbol{z}^{(k)})
abla_{\phi} \log q_{\phi}(oldsymbol{z}^{(k)})$$

イロト イ理ト イヨト イヨト 三国

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

• • = • • = •

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 (Maddison et al., 2017; Jang et al., 2017)
- Variance reduction by adding control variates (a.k.a. baselines)

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

where $\mu_c = \nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z)}[c(z)] = \mathbb{E}_{z \sim q_{\phi}(z)}[c(z)\nabla_{\phi} \log q_{\phi}(z)]$ is known

- REBAR (Tucker et al., 2017)
- RELAX (Grathwohl et al., 2018)

Gradient estimation for discrete latent variables

Our ideas:

- Spike gradient:
 - Don't move unless you are pretty sure which direction to move
 - If you do move, move with large spikes
 - The temporal average of the spikes shall code the temporal evolution of the true gradient
- Variance reduction by sharing common random numbers between different expectations
- No need to construct baselines (control variates), as the function *f* itself will be used to construct a baseline
- Augmentation + REINFORCE + merge is how we derive such a gradient estimator

通 ト イヨ ト イヨト

Exponential, Gumbel, and categorical random variables

• If $x_i \sim \text{Exp}(\lambda_i)$ are independent exponential random variables for i = 1, ..., M, then

$$m{P}ig(i = {
m arg\,min}_j \, x_jig) = m{P}ig(x_i < x_j, \; orall \, j
eq iig) = \lambda_i \Big/{\sum_{i=1}^M \lambda_i} \;\;.$$

• Note $x \sim \text{Exp}(\lambda)$ can be reparameterized as

$$x = \epsilon / \lambda, \ \epsilon \sim \mathsf{Exp}(1),$$

where $\epsilon \sim \text{Exp}(1)$ can be equivalently generated as

$$\epsilon = - \log u, \,\, u \sim {\sf Uniform}(0,1)$$

・ロン ・四 ・ ・ ヨン ・ ヨン

Exponential, Gumbel, and categorical random variables

• If

 $x_i \sim \operatorname{Exp}(\lambda_i),$

then we have

$$\begin{aligned} \arg\min_{i} x_{i} \stackrel{d}{=} \arg\min_{i} \{-\log u_{i}/\lambda_{i}\}, & u_{i} \stackrel{iid}{\sim} \mathsf{Uniform}(0,1) \\ &= \arg\max_{i} \{\log \lambda_{i} - \log(-\log u_{i})\} \\ &\stackrel{d}{=} \arg\max_{i} \{\log \lambda_{i} + \epsilon_{i}\}, \quad \epsilon_{i} \sim \mathsf{Gumbel}(0,1) \end{aligned}$$

 Note if u ~ Uniform(0, 1), then e = − log(− log u) follows the Gumbel(0, 1) distribution (Type-I extreme-value distribution)

(日) (同) (三) (三)

Augmentation of categorical random variable

Denoting
$$\sigma(\phi) = \left(\frac{e^{\phi_1}}{\sum_{m=1}^M e^{\phi_m}}, \dots, \frac{e^{\phi_M}}{\sum_{m=1}^M e^{\phi_m}}\right)$$
, categorical $z \sim \text{Discrete}(\sigma(\phi))$ can be augmented as
 $z = \arg\min_{i \in \{1,\dots,M\}} e_i$, where $e_i \sim \text{Exp}(e^{\phi_i})$

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

$$\begin{split} \mathcal{E}(\phi) &= \mathbb{E}_{z \sim \mathsf{Discrete}(\sigma(\phi))}[f(z)] \\ &= \mathbb{E}_{e_1 \sim \mathsf{Exp}(e^{\phi_1}), \dots, e_M \sim \mathsf{Exp}(e^{\phi_M})}[f(\arg\min_i e_i)] \end{split}$$

Mingyuan Zhou (UT-McCombs)

July 2018 10 / 30

(日) (周) (三) (三)

Use REINFORCE we have

$$egin{aligned}
abla_{\phi_m} \mathcal{E}(\phi) &= \mathbb{E}_{e_1 \sim \mathsf{Exp}(e^{\phi_1}), ..., e_M \sim \mathsf{Exp}(e^{\phi_M})}[f(rg\min_i e_i)
abla_{\phi_m} \log \mathsf{Exp}(e_m; e^{\phi_m})] \ &= \mathbb{E}_{e_1 \sim \mathsf{Exp}(e^{\phi_1}), ..., e_M \sim \mathsf{Exp}(e^{\phi_M})}[f(rg\min_i e_i)(1 - e_m e^{\phi_m})] \end{aligned}$$

Since the exponential random variable $x \sim \text{Exp}(e^{\phi})$ can be reparameterized as $x = \epsilon e^{-\phi}$, $\epsilon \sim \text{Exp}(1)$, we have

$$\nabla_{\phi_m} \mathcal{E}(\phi) = \mathbb{E}_{\substack{\epsilon_1, \dots, \epsilon_M \stackrel{iid}{\sim} \mathsf{Exp}(1)}} [f(\arg\min_i \epsilon_i e^{-\phi_i})(1-\epsilon_m)]$$

.

REINFORCE in augmented space

A key observation is that we may choose M as reference category and rewrite the gradient estimator

$$\nabla_{\phi_m} \mathcal{E}(\phi) = \mathbb{E}_{\substack{\epsilon_1, \dots, \epsilon_M \stackrel{\text{iid}}{\sim} \mathsf{Exp}(1)}} [f(\arg\min_i \epsilon_i e^{-\phi_i})(1-\epsilon_m)]$$

as

$$\nabla_{\phi_m} \mathcal{E}(\phi) = \mathbb{E}_{\substack{\epsilon_1, \dots, \epsilon_M \stackrel{iid}{\sim} \mathsf{Exp}(1)}} [f(\arg\min_i \epsilon_{(m \leftrightarrows M)_i} e^{-\phi_i})(1 - \epsilon_M)]$$

where (m = M) denotes a vector of indices constructed by swapping the *m*-th and *M*-th elements of vector $(1, \ldots, M)$, which means

$$\begin{cases} (m \leftrightarrows M)_M = m \\ (m \leftrightarrows M)_m = M \\ (m \leftrightarrows M)_i = i, & \text{if } i \notin \{m, M\} \end{cases}$$

Merge the gradients

As $\sigma(\phi - \phi_M \mathbf{1}_M) = \sigma(\phi)$, one may update $\tilde{\phi}_m = \phi_m - \phi_M$ for $m \le M - 1$ and set $\tilde{\phi}_M = 0$. Denoting $\tilde{\phi} = (\tilde{\phi}_1, \dots, \tilde{\phi}_{M-1})' = \mathbf{A}\phi$, where $\mathbf{A} = [\operatorname{diag}(\mathbf{1}_{M-1}), -\mathbf{1}_{M-1}]$, we have

$$\begin{split} \nabla_{\phi} \mathcal{E}(\phi)' = &\nabla_{\tilde{\phi}} \mathcal{E}([\tilde{\phi}',0]')' \frac{\partial \tilde{\phi}}{\partial \phi} = \nabla_{\tilde{\phi}} \mathcal{E}([\tilde{\phi}',0]')' \mathbf{A} \\ \nabla_{\tilde{\phi}_m} \mathcal{E}([\tilde{\phi}',0]') = &\frac{1}{M} \sum_{j=1}^M (\nabla_{\phi_m} \mathcal{E}(\phi) - \nabla_{\phi_j} \mathcal{E}(\phi)) \\ = &\nabla_{\phi_m} \mathcal{E}(\phi) - \frac{1}{M} \sum_{j=1}^M \nabla_{\phi_j} \mathcal{E}(\phi). \end{split}$$

Notice previously we have

$$\nabla_{\phi_m} \mathcal{E}(\phi) = \mathbb{E}_{\substack{\epsilon_1, \dots, \epsilon_M \stackrel{iid}{\sim} \mathsf{Exp}(1)}} [f(\arg\min_i \epsilon_{(m \leftrightarrows M)_i} e^{-\phi_i})(1 - \epsilon_M)].$$

Mingyuan Zhou (UT-McCombs)

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ののの

Merge the gradients

Combining the two equations the ARM estimator is

$$abla_{ ilde{\phi}_m} \mathcal{E}([ilde{\phi}', 0]') = \mathbb{E}_{\epsilon_1, ..., \epsilon_M} igar_{\sim}^{iid} \mathsf{Exp}(1)} [f_\Delta(\epsilon, \phi, m)(1-\epsilon_M)],$$

Common random numbers are shared to compute:

$$f_{\Delta}(\epsilon, \phi, m) = f(\arg\min_{i} \epsilon_{(m \leftrightarrows M)_{i}} e^{-\phi_{i}}) - \frac{1}{M} \sum_{j=1}^{M} f(\arg\min_{i} \epsilon_{(j \oiint M)_{i}} e^{-\phi_{i}})$$
$$= \frac{1}{M} \sum_{j \neq m} \left[f(\arg\min_{i} \epsilon_{(m \leftrightarrows M)_{i}} e^{-\phi_{i}}) - f(\arg\min_{i} \epsilon_{(j \oiint M)_{i}} e^{-\phi_{i}}) \right]$$

Note that

$$\mathbb{E}_{\epsilon_1,...,\epsilon_M \stackrel{iid}{\sim} \mathsf{Exp}(1)}[f_\Delta(\epsilon,\phi,m)] = 0,$$

thus there is no need to add control variates

< ロ > < 同 > < 三 > < 三

Merge the gradients

• $\epsilon_1,\ldots,\epsilon_M \stackrel{\textit{iid}}{\sim} \mathsf{Exp}(1)$ is the same in distribution as

 $\epsilon_i = \pi_i \epsilon$, for $i = 1, \dots, M$, where $\pi \sim \text{Dirichlet}(\mathbf{1}_M), \ \epsilon \sim \text{Gamma}(M, 1)$,

•
$$\arg\min_i \pi_{(m \subseteq M)_i} e^{-\phi_i} = \arg\min_i \epsilon \pi_{(m \subseteq M)_i} e^{-\phi_i}$$

• The gradient can be re-expressed as

$$abla_{ ilde{\phi}_m} \mathcal{E}(\phi) =
abla_{ ilde{\phi}_m} \mathcal{E}([ilde{\phi}', 0]') = \mathbb{E}_{m{\pi} \sim \mathsf{Dirichlet}(\mathbf{1}_M)}[f_\Delta(m{\pi}, \phi, m)(1 - M \pi_M)]$$

where

$$f_{\Delta}(\pi,\phi,m) = f\left(\underset{i\in\{1,\dots,M\}}{\operatorname{arg\,min}} \pi_{(m\leftrightarrows M)_i}e^{-\phi_i}\right) - \frac{1}{M}\sum_{j=1}^M f\left(\underset{i\in\{1,\dots,M\}}{\operatorname{arg\,min}} \pi_{(j\leftrightarrows M)_i}e^{-\phi_i}\right) \\ = \frac{1}{M}\sum_{j\neq m} \left[f\left(\underset{i\in\{1,\dots,M\}}{\operatorname{arg\,min}} \pi_{(m\leftrightarrows M)_i}e^{-\phi_i}\right) - f\left(\underset{i\in\{1,\dots,M\}}{\operatorname{arg\,min}} \pi_{(j\leftrightharpoons M)_i}e^{-\phi_i}\right)\right]$$

Mingyuan Zhou (UT-McCombs)

(日) (周) (三) (三)

ARM gradient for binary random variable

For a binary random variable, the gradient of $\mathbb{E}_{z \sim \text{Bernoulli}(\sigma(\phi))}[f(z)]$ with respect to ϕ can be expressed as

$$\nabla_{\phi} \mathcal{E}(\phi) = \mathbb{E}_{u \sim \text{Uniform}(0,1)} [f_{\Delta}(u,\phi)(u-1/2)]$$

$$f_{\Delta}(u,\phi) = f(\mathbf{1}[u > \sigma(-\phi)]) - f(\mathbf{1}[u < \sigma(\phi)]).$$

Proposition

(i) $f_{\Delta}(u, \phi) = 0$ with probability $\sigma(|\phi|) - \sigma(-|\phi|)$, $f_{\Delta}(u, \phi) = f(1) - f(0)$ with probability $1 - \sigma(|\phi|)$, and $f_{\Delta}(u, \phi) = f(0) - f(1)$ with probability $\sigma(-|\phi|)$.

(ii)
$$g_{ARM}(u,\phi) = f_{\Delta}(u,\phi)(u-1/2)$$
 is unbiased with
 $\mathbb{E}_{u \sim Uniform(0,1)}[g_{ARM}(u,\phi)] = \mathbb{E}_{z \sim Bernoulli(\sigma(\phi))}[g_{REINFORCE}(z,\phi)].$

(iii) $g_{ARM}(u, \phi)$ reaches its largest variance at 0.039788 $[f(1) - f(0)]^2$ when $\frac{P(f_\Delta = 0)}{P(f_\Delta \neq 0)}$ is equal to the golden ratio $\frac{\sqrt{5}+1}{2}$.

(iv)
$$\frac{\sup_{\phi} \operatorname{Var}[g_{ARM}]}{\sup_{\phi} \operatorname{Var}[g_{REINFORCE}]} \leq \frac{16}{25} (1 - 2\frac{f(0)}{f(0) + f(1)})^2 \text{ and } \sup_{\phi} \operatorname{Var}[g_{ARM}] \leq \frac{1}{25} [f(1) - f(0)]^2.$$

A simple example

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



Mingyuan Zhou (UT-McCombs)

 $ARM-\nabla$

July 2018 17 / 30

3

(日) (周) (三) (三)

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.

 • A latent variable model with multiple stochastic hidden layers is

$$oldsymbol{x} \sim p_{oldsymbol{ heta}_0}(oldsymbol{x} \mid oldsymbol{b}_1), \ oldsymbol{b}_1 \sim p_{oldsymbol{ heta}_1}(oldsymbol{b}_1 \mid oldsymbol{b}_2), \dots, oldsymbol{b}_{\mathcal{T}} \sim p_{oldsymbol{ heta}_{\mathcal{T}}}(oldsymbol{b}_{\mathcal{T}}),$$

• The joint likelihood is

$$p(\boldsymbol{x}, \boldsymbol{b}_{1:T} \mid \boldsymbol{\theta}_{0:T}) = p_{\boldsymbol{\theta}_0}(\boldsymbol{x} \mid \boldsymbol{b}_1) \Big[\prod_{t=1}^{T-1} p_{\boldsymbol{\theta}_t}(\boldsymbol{b}_t \mid \boldsymbol{b}_{t+1}) \Big] p_{\boldsymbol{\theta}_T}(\boldsymbol{b}_T).$$

→ Ξ →

VAE with multiple discrete stochastic layers

• The encoder is designed as

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

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

$$\mathcal{E}(\boldsymbol{w}_{1:T}) = \mathbb{E}_{\boldsymbol{b}_{1:T} \sim q_{\boldsymbol{w}_{1:T}}(\boldsymbol{b}_{1:T} \mid \boldsymbol{x})} [f(\boldsymbol{b}_{1:T})], \text{ where}$$

$$f(\boldsymbol{b}_{1:T}) = \log p_{\boldsymbol{\theta}_0}(\boldsymbol{x} \mid \boldsymbol{b}_1) + \log p_{\boldsymbol{\theta}_{1:T}}(\boldsymbol{b}_{1:T}) - \log q_{\boldsymbol{w}_{1:T}}(\boldsymbol{b}_{1:T} \mid \boldsymbol{x}).$$

First, to compute the gradient with respect to \boldsymbol{w}_1 , since

$$\mathcal{E}(\boldsymbol{w}_{1:T}) = \mathbb{E}_{q(\boldsymbol{b}_1)} \mathbb{E}_{q(\boldsymbol{b}_{2:T} \mid \boldsymbol{b}_1)}[f(\boldsymbol{b}_{1:T})]$$

we have

$$\nabla_{\boldsymbol{w}_1} \mathcal{E}(\boldsymbol{w}_{1:T}) = \mathbb{E}_{\boldsymbol{u}_1 \sim \text{Uniform}(0,1)}[f_{\Delta}(\boldsymbol{u}_1, \mathcal{T}_{\boldsymbol{w}_1}(\boldsymbol{x}))(\boldsymbol{u}_1 - 1/2)] \nabla_{\boldsymbol{w}_1} \mathcal{T}_{\boldsymbol{w}_1}(\boldsymbol{x}),$$

where

$$f_{\Delta}(\boldsymbol{u}_{1}, \mathcal{T}_{\boldsymbol{w}_{1}}(\boldsymbol{x})) = \mathbb{E}_{\boldsymbol{b}_{2:T} \sim q(\boldsymbol{b}_{2:T} \mid \boldsymbol{b}_{1}), \ \boldsymbol{b}_{1} = \mathbf{1}[\boldsymbol{u}_{1} > \sigma(-\mathcal{T}_{\boldsymbol{w}_{1}}(\boldsymbol{x}))])[f(\boldsymbol{b}_{1:T})]} \\ - \mathbb{E}_{\boldsymbol{b}_{2:T} \sim q(\boldsymbol{b}_{2:T} \mid \boldsymbol{b}_{1}), \ \boldsymbol{b}_{1} = \mathbf{1}[\boldsymbol{u}_{1} < \sigma(\mathcal{T}_{\boldsymbol{w}_{1}}(\boldsymbol{x}))])[f(\boldsymbol{b}_{1:T})]}$$

< ロ > < 同 > < 三 > < 三

ARM gradient of the ELBO

Second, to compute the gradient with respect to \boldsymbol{w}_t , where $2 \le t \le T - 1$, since

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

we have

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

where

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

(日) (周) (三) (三)

Finally, to compute the gradient with respect to $\boldsymbol{w}_{\mathcal{T}}$, we have

$$\nabla_{\boldsymbol{w}_{\mathcal{T}}} \mathcal{E}(\boldsymbol{w}_{1:\mathcal{T}})$$

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

$$f_{\Delta}(\boldsymbol{u}_{T}, \mathcal{T}_{\boldsymbol{w}_{T}}(\boldsymbol{b}_{T-1}), \boldsymbol{b}_{1:T-1}) = f(\boldsymbol{b}_{1:T-1}, \boldsymbol{b}_{T} = \mathbf{1}[\boldsymbol{u}_{T} > \sigma(-\mathcal{T}_{\boldsymbol{w}_{T}}(\boldsymbol{b}_{T-1}))])) \\ - f(\boldsymbol{b}_{1:T-1}, \boldsymbol{b}_{T} = \mathbf{1}[\boldsymbol{u}_{T} < \sigma(\mathcal{T}_{\boldsymbol{w}_{T}}(\boldsymbol{b}_{T-1}))])$$

< ロ > < 同 > < 三 > < 三



Figure 2: Test negative ELBOs on MNIST 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.

(日) (同) (三) (三)

			ARM	RELAX	REBAR	ST Gumbel-Softmax
Bernoulli	Nonlinear	MNIST OMNIGLOT	101.3 129.5	110.9 128.2	111.6 128.3	112.5 140.7
	Linear	MNIST OMNIGLOT	110.3 124.2	122.1 124.4	123.2 124.9	129.2 129.8
	Two layers	MNIST OMNIGLOT	98.2 118.3	114.0 119.1	113.7 118.8	NA NA
Categorical	Nonlinear	MNIST OMNIGLOT	105.8 121.9	NA NA	NA NA	107.9 127.6

Table 2: Test negative ELBOs of discrete VAEs trained with four different stochastic gradient estimators.

イロト イポト イヨト イヨト 二日



Figure: Comparison of the negative ELBOs for categorical variational auto-encoders trained by ARM and ST Gumbel-softmax on MNIST and OMNIGLOT, using the "Nonlinear" network.

Gradient estimator	ARM	ST	1/2	Annealed ST	ST Gumbel-S.	SF	MuProp
$-\log p(\boldsymbol{x}_{l} \mid \boldsymbol{x}_{u})$	54.8	56.1	57.2	58.7	59.3	72.0	56.7

→ 3 → 4 3

Image: A matrix

MLE with multiple discrete stochastic layers

• The log marginal likelihood can be expressed as

$$\begin{split} \log p_{\boldsymbol{\theta}_{0:\mathcal{T}}}(\boldsymbol{x}) &= \log \mathbb{E}_{\boldsymbol{b}_{1:\mathcal{T}} \sim p_{\boldsymbol{\theta}_{1:\mathcal{T}}}(\boldsymbol{b}_{1:\mathcal{T}})} [p_{\boldsymbol{\theta}_{0}}(\boldsymbol{x} \mid \boldsymbol{b}_{1})] \\ &\geq \mathcal{E}(\boldsymbol{\theta}_{1:\mathcal{T}}) = \mathbb{E}_{\boldsymbol{b}_{1:\mathcal{T}} \sim p_{\boldsymbol{\theta}_{1:\mathcal{T}}}(\boldsymbol{b}_{1:\mathcal{T}})} [\log p_{\boldsymbol{\theta}_{0}}(\boldsymbol{x} \mid \boldsymbol{b}_{1})]. \end{split}$$

• For stochastic binary network

$$p_{\boldsymbol{\theta}_t}(\boldsymbol{b}_t \mid \boldsymbol{b}_{t+1}) = \text{Bernoulli}(\boldsymbol{b}_t; \sigma(\mathcal{T}_{\boldsymbol{\theta}_t}(\boldsymbol{b}_{t+1}))),$$

• The gradient of the lower bound can be expressed as

$$\nabla_{\boldsymbol{\theta}_{t}} \mathcal{E}(\boldsymbol{\theta}_{1:T}) = \mathbb{E}_{\boldsymbol{\rho}(\boldsymbol{b}_{t+1:T})} \left[\mathbb{E}_{\boldsymbol{u}_{t}} \left[f_{\Delta}(\boldsymbol{u}_{t}, \mathcal{T}_{\boldsymbol{\theta}_{t}}(\boldsymbol{b}_{t+1}), \boldsymbol{b}_{t+1:T})(\boldsymbol{u}_{t} - 1/2) \right] \nabla_{\boldsymbol{\theta}_{t}} \mathcal{T}_{\boldsymbol{\theta}_{t}}(\boldsymbol{b}_{t+1}) \right],$$

$$f_{\Delta}(\boldsymbol{u}_{t}, \mathcal{T}_{\boldsymbol{\theta}_{t}}(\boldsymbol{b}_{t+1}), \boldsymbol{b}_{t+1:T}) = \mathbb{E}_{\boldsymbol{b}_{1:t-1} \sim \boldsymbol{\rho}(\boldsymbol{b}_{1:t-1} \mid \boldsymbol{b}_{t}), \ \boldsymbol{b}_{t} = \mathbf{1}[\boldsymbol{u}_{t} > \boldsymbol{\sigma}(-\mathcal{T}_{\boldsymbol{\theta}_{t}}(\boldsymbol{b}_{t+1}))]) \left[\log \boldsymbol{p}_{\boldsymbol{\theta}_{0}}(\boldsymbol{x} \mid \boldsymbol{b}_{1}) \right]$$

$$- \mathbb{E}_{\boldsymbol{b}_{1:t-1} \sim \boldsymbol{\rho}(\boldsymbol{b}_{1:t-1} \mid \boldsymbol{b}_{t}), \ \boldsymbol{b}_{t} = \mathbf{1}[\boldsymbol{u}_{t} < \boldsymbol{\sigma}(\mathcal{T}_{\boldsymbol{\theta}_{t}}(\boldsymbol{b}_{t+1}))]) \left[\log \boldsymbol{p}_{\boldsymbol{\theta}_{0}}(\boldsymbol{x} \mid \boldsymbol{b}_{1}) \right].$$

Mingyuan Zhou (UT-McCombs)

MLE with multiple discrete stochastic layers

- Predict lower half of MNIST digit x₁ given the upper half x_u;
- Maximizing the conditional likelihood $p_{\theta_{0:2}}(\mathbf{x}_{I} | \mathbf{x}_{u})$
- Approximate log $p_{\boldsymbol{\theta}_{0:2}}(\boldsymbol{x}_{l} \mid \boldsymbol{x}_{u})$ with

$$\log \frac{1}{K} \sum_{k=1}^{K} \mathsf{Bernoulli}(\boldsymbol{x}_{l}; \sigma(\mathcal{T}_{\boldsymbol{\theta}_{0}}(\boldsymbol{b}_{1}^{(k)})))$$

where $\boldsymbol{b}_1^{(k)} \sim \text{Bernoulli}(\sigma(\mathcal{T}_{\boldsymbol{\theta}_1}(\boldsymbol{b}_2^{(k)}))), \ \boldsymbol{b}_2^{(k)} \sim \text{Bernoulli}(\sigma(\mathcal{T}_{\boldsymbol{\theta}_2}(\boldsymbol{x}_u))).$

• Training with K=1 and on the test approximate negative log-likelihood $-\log p_{\theta_{0:2}}(\mathbf{x}_{I} | \mathbf{x}_{u})$ with K = 1000

MLE with multiple discrete stochastic layers



Table 3: For the MNIST conditional distribution estimation benchmark task, comparison of the test negative log-likelihood between various gradient estimators, with the best results in [14, 20] reported here.

Gradient estimator	ARM	ST	1/2	Annealed ST	ST Gumbel-S.	SF	MuProp
$-\log p(\boldsymbol{x}_l \boldsymbol{x}_u)$	54.8	56.1	57.2	58.7	59.3	72.0	56.7

Mingyuan Zhou (UT-McCombs)

July 2018

29 / 30

Thank you!

・ロト ・回ト ・ヨト ・ヨ