

Modern Computational Statistics

Lecture 16: Advanced VI



Cheng Zhang

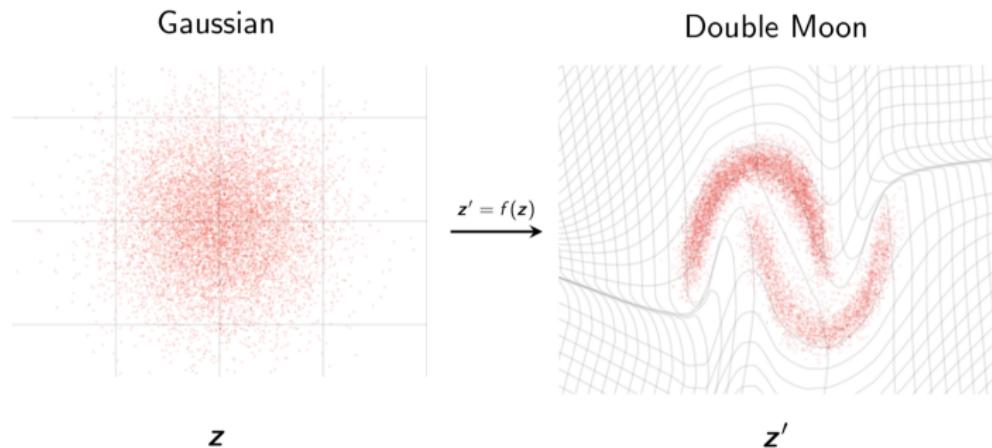
School of Mathematical Sciences, Peking University

November 20, 2019

- ▶ The approximation accuracy of VI depends on the **expressive power** of the approximating distributions.
- ▶ Ideally, we want a rich variational family of distributions that provide accurate approximation while maintaining the computational efficiency and scalability.
- ▶ In this lecture, we will discuss some recent techniques for improving the flexibility of variational approximations.
- ▶ We will also talk about how to combine the two engines of modern Bayesian inference, MCMC and VI, to get the best of both worlds.

- ▶ VI requires the approximating distributions to have the following properties
 - ▶ Analytic density
 - ▶ Easy to sample
- ▶ Many simple distributions satisfy the above properties, e.g., Gaussian, general exponential family distributions. Therefore, they are commonly used in VI.
- ▶ Unfortunately, the posterior distribution could be much more complex (highly skewed, multi-modal, etc).
- ▶ **How can we improve the complexity of our variational approximations while maintaining the desired properties?**

- Idea: Map simple distributions to complex distributions via learnable transforms.



Change of Variables

Assume that the mapping between z and x , given by

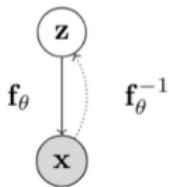
$f : \mathbb{R}^n \rightarrow \mathbb{R}^n$, is invertible such that $x = f(z)$ and $z = f^{-1}(x)$

$$p_x(x) = p_z(f^{-1}(x)) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|$$

- ▶ x, z need to be continuous and have the same dimension.
For example, if $x \in \mathbb{R}^n$ then $z \in \mathbb{R}^n$
- ▶ For any invertible matrix A , $\det(A^{-1}) = \det(A)^{-1}$

$$p_x(x) = p_z(z) \left| \det \left(\frac{\partial f(z)}{\partial z} \right) \right|^{-1}$$

- ▶ Consider a directed, latent-variable model over observed variables x and latent variables z .
- ▶ In a normalizing flow model, the mapping between z and x , given by $f_\theta : \mathbb{R}^n \mapsto \mathbb{R}^n$, is deterministic and invertible such that $x = f_\theta(z)$ and $z = f_\theta^{-1}(x)$



- ▶ Using change of variables, the probability $p(x)$ is given by

$$p_x(x|\theta) = p_z(z) \left| \det \left(\frac{\partial f_\theta(z)}{\partial z} \right) \right|^{-1}$$

- ▶ **Normalizing Transforms:** Change of variables gives a normalized density after applying an invertible transformation
- ▶ **Flow:** Invertible transformations can be composed with each other

$$z_k = f_k(z_{k-1}), \quad k = 1, \dots, K$$

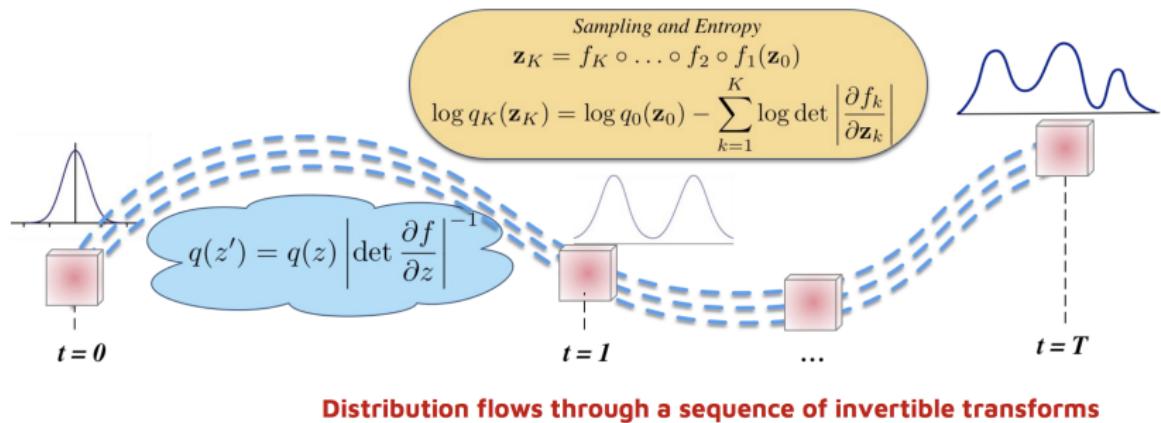
- ▶ The log-likelihood of z_K

$$\log p_K(z_K) = \log p_0(z_0) - \sum_{k=1}^K \log \left| \det \left(\frac{\partial f_k(z_{k-1})}{z_{k-1}} \right) \right|$$

Remark: for simplicity, we omit the parameters for each of these transformations f_1, f_2, \dots, f_K .

Exploit the rule for change of variables

- ▶ Start with a simple distribution for z_0 (e.g., Gaussian).
- ▶ Apply a sequence of K invertible transformations.



Adapted from Mohamed and Rezende, 2017

- ▶ Planar flow (Rezende and Mohamed, 2015).

$$x = f_\theta(z) = z + uh(w^\top z + b)$$

parameterized by $\theta = (w, u, b)$ where h is a non-linear function

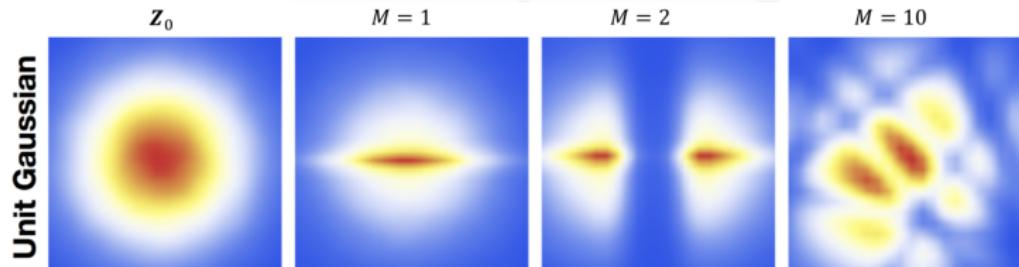
- ▶ Absolute value of the determinant of the Jacobian

$$\begin{aligned} \left| \det \frac{\partial f_\theta(z)}{\partial z} \right| &= \left| \det(I + h'(w^\top z + b)uw^\top) \right| \\ &= \left| 1 + h'(w^\top z + b)u^\top w \right| \end{aligned}$$

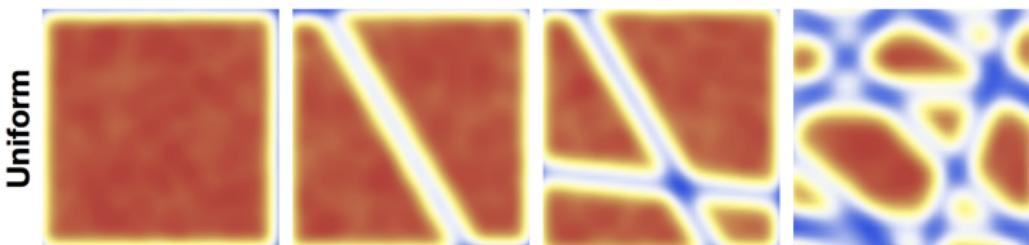
- ▶ Need to restrict parameters and non-linearity for the mapping to be invertible. For example,

$$h(\cdot) = \tanh(\cdot), \quad h'(w^\top z + b)u^\top w \geq 1$$

- Base distribution: Gaussian



- Base distribution: Uniform



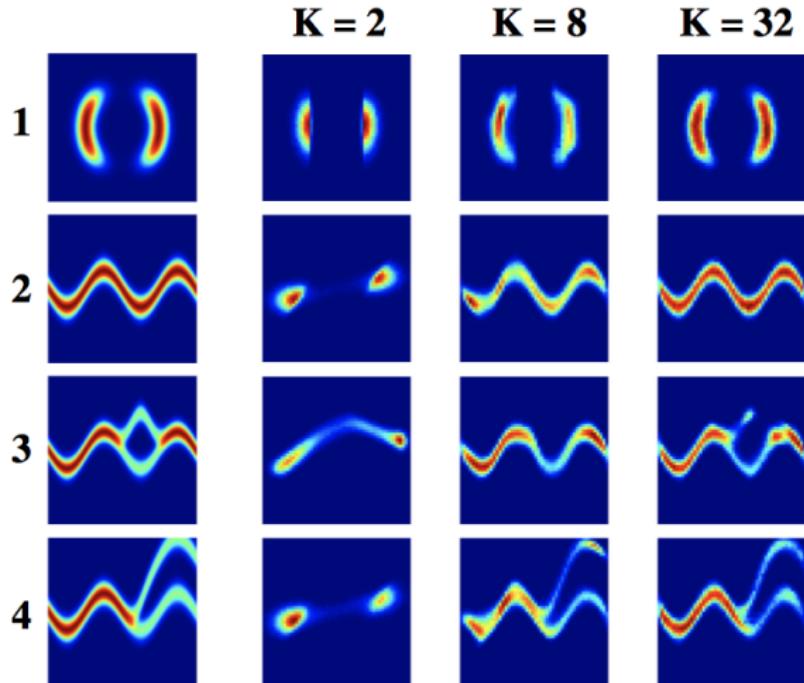
- 10 planar transformations can transform simple distributions into a more complicated one.

- ▶ Learning via maximizing the ELBO

$$\begin{aligned} L &= \mathbb{E}_{q_K(z_K)} \log \frac{p(x, z_K)}{q_K(z_K)} \\ &= \mathbb{E}_{q_0(z_0)} \log p(x, z_K) - \mathbb{E}_{q_0(z_0)} \log q_0(z_0) \\ &\quad - \sum_{k=1}^K \mathbb{E}_{q_0(z_0)} \log \left| \det \left(\frac{\partial f_k(z_{k-1})}{\partial z_{k-1}} \right) \right| \end{aligned}$$

- ▶ Exact likelihood evaluation via inverse transformation and change of variable formula
- ▶ Sampling via forward transformation

$$z_0 \sim q_0(z_0), \quad z_K = f_K \circ f_{K-1} \circ \cdots \circ f_1(z_0)$$



Adapted from Rezenda and Mohamed, 2015

- ▶ Simple initial distribution $q_0(z_0)$ that allows for efficient sampling and tractable likelihood evaluation, e.g., Gaussian
- ▶ Sampling requires efficient evaluation of

$$z_k = f_k(z_{k-1}), \quad k = 1, \dots, K$$

- ▶ Likelihood computation also requires the evaluation of determinants of $n \times n$ Jacobian matrices $\sim \mathcal{O}(n^3)$, prohibitively expensive within a learning loop!
- ▶ Design transformations so that the resulting Jacobian matrix has special structure. For example
 - ▶ lower rank update to identity as in planar flows.
 - ▶ triangular matrix whose determinant is just the product of the diagonal entries, i.e., an $\mathcal{O}(n)$ operation.

- ▶ NICE or Nonlinear Independent Components Estimation (Dinh et al., 2014) composes two kinds of invertible transformations: additive coupling layers and rescaling layers
- ▶ Real-NVP (Dinh et al., 2017)
- ▶ Inverse Autoregressive Flow (Kingma et al., 2016)
- ▶ Masked Autoregressive Flow (Papamakarios et al., 2017)

- ▶ Partition the variable z into two disjoint subsets

$$z = z_{1:d} \cup z_{d+1:n}$$

- ▶ Forward mapping $z \mapsto x$:

$$x_{1:d} = z_{1:d}, \quad x_{d+1:n} = z_{d+1:n} + m_\theta(z_{1:d})$$

where $m_\theta : \mathbb{R}^d \mapsto \mathbb{R}^{n-d}$ is a neural network with parameters θ

- ▶ Backward mapping $x \mapsto z$:

$$z_{1:d} = x_{1:d}, \quad z_{d+1:n} = x_{d+1:n} - m_\theta(x_{1:d})$$

- ▶ Forward/Backward mapping is **volume preserving**: the determinant of the Jacobian is 1.

- ▶ Additive coupling layers are composed together (with arbitrary partitions of variables in each layer)
- ▶ Final layer of NICE uses a rescaling transformation
- ▶ Forward mapping $z \mapsto x$:

$$x_i = s_i z_i, \quad i = 1, \dots, n$$

where $s_i > 0$ is the scaling factor for the i -th dimension.

- ▶ Backward mapping $x \mapsto z$:

$$z_i = \frac{x_i}{s_i}, \quad i = 1, \dots, n$$

- ▶ Jacobian of forward mapping:

$$J = \text{diag}(s), \quad \det(J) = \prod_{i=1}^n s_i.$$

- ▶ Forward mapping $z \mapsto x$:

$$x_{1:d} = z_{1:d}, \quad x_{d+1:n} = z_{d+1:n} \odot \exp(\alpha_\theta(z_{1:d})) + \mu_\theta(z_{1:d})$$

where α_θ and μ_θ are both neural networks.

- ▶ Backward mapping $x \mapsto z$:

$$z_{1:d} = x_{1:d}, \quad z_{d+1:n} = \exp(-\alpha_\theta(x_{1:d})) \odot (x_{d+1:n} - \mu_\theta(x_{1:d}))$$

- ▶ The determinant of the Jacobian of forward mapping

$$\det \left(\frac{\partial x}{\partial z} \right) = \exp \left(\sum \alpha_\theta(z_{1:d}) \right)$$

- ▶ Non-volume preserving transformation in general since determinant can be less than or greater than 1.

- ▶ Consider a Gaussian autoregressive model

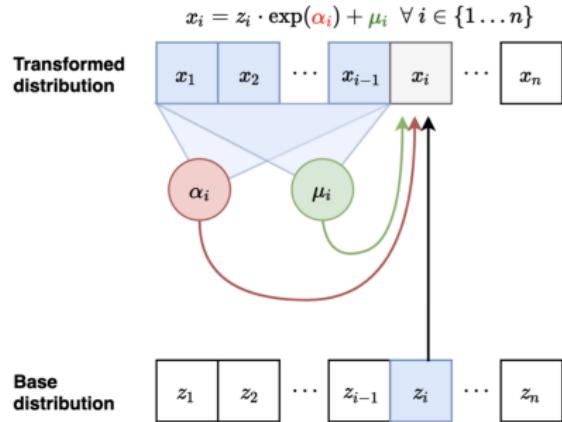
$$p(x) = \prod_{i=1}^n p(x_i | x_{<i})$$

where $p(x_i | x_{<i}) = \mathcal{N}(\mu_i(x_{1:i-1}), \exp(\alpha_i(x_{1:i-1}))^2)$. μ_i and α_i are neural networks for $i > 1$ and constants for $i = 1$.

- ▶ Sequential sampling:

$$z_i \sim \mathcal{N}(0, 1), \quad x_i = \exp(\alpha_i(x_{1:i-1})) z_i + \mu_i(x_{1:i-1}), \quad i = 1, \dots, n$$

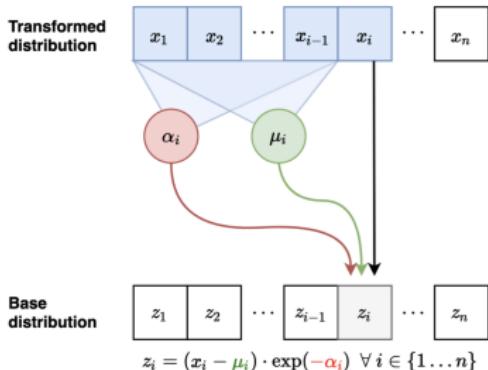
- ▶ **Flow interpretation:** transforms samples from the standard Gaussian to those generated from the model via invertible transformations (parameterized by μ_i, α_i)



- ▶ Forward mapping from $z \mapsto x$:

$$x_i = \exp(\alpha_i(x_{1:i-1}))z_i + \mu_i(x_{1:i-1}), \quad i = 1, \dots, n$$

- ▶ Like autoregressive models, sampling is sequential and slow ($\mathcal{O}(n)$)

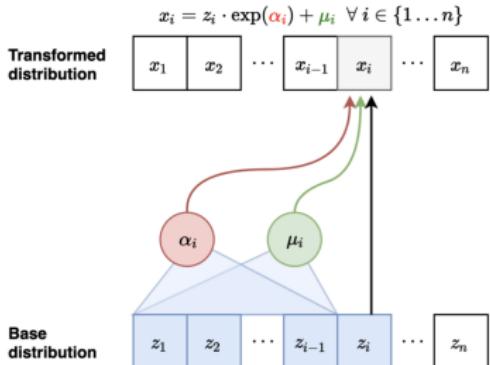


- Inverse mapping from $x \mapsto z$: shift and scale

$$z_i = (x_i - \mu_i(x_{1:i-1})) / \exp(\alpha_i(x_{1:i-1})), \quad i = 1, \dots, n$$

Note that this can be done in parallel.

- Jacobian is lower diagonal, hence determinant can be computed efficiently.
- Likelihood evaluation is easy and parallelizable.



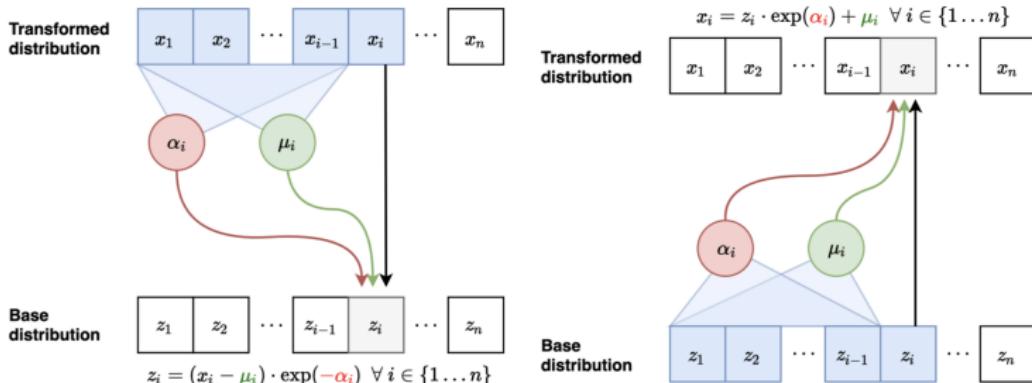
- ▶ Forward mapping from $z \mapsto x$ (parallel):

$$x_i = \exp(\alpha_i(z_{1:i-1}))z_i + \mu_i(z_{1:i-1}), \quad i = 1, \dots, n$$

- ▶ Backward mapping from $x \mapsto z$ (sequential):

$$z_i = (x_i - \mu_i(z_{1:i-1}))/\exp(\alpha_i(z_{1:i-1}))$$

- ▶ Fast to sample from, slow to evaluate likelihoods of data points. However, likelihood evaluation for a sampled point is fast.



Inverse pass of MAF (left) vs. Forward pass of IAF (right)

- ▶ Interchanging z and x in the inverse transformation of MAF gives the forward transformation of IAF.
- ▶ Similarly, forward transformation of MAF is inverse transformation of IAF.

- ▶ Transform simple distributions into more complex distributions via change of variables
- ▶ Jacobian of transformations should have tractable determinant for efficient learning and density estimation
- ▶ Computational tradeoff in evaluating forward and inverse transformations
 - ▶ **MAF**: Fast likelihood evaluation, slow sampling, more suited for MLE based training, density estimation.
 - ▶ **IAF**: Fast sampling, slow likelihood evaluation, more suited for variational inference, real time generation.
 - ▶ **NICE** and **RealNVP**: Fast on both side, but generally less flexible than the others.

- MCMC approximates the posterior through a sequence of transitions

$$z_0 \sim q(z_0), \quad z_t \sim q(z_t | z_{t-1}, x), \quad t = 1, 2, \dots$$

where the transition kernel satisfies the detailed balance condition

$$p(x, z_{t-1})q(z_t | z_{t-1}, x) = p(x, z_t)q(z_{t-1} | z_t, x)$$

- Pros
 - automatically adapts to true posterior
 - asymptotically unbiased
- Cons
 - slow convergence, hard to assess quality
 - tuning headaches

- ▶ Each iteration in MCMC can be viewed as a mapping $z_{t-1} \mapsto z_t$, and the marginal likelihood of z_T is

$$q(z_T|x) = \int q(z_0|x) \prod_{t=1}^T q(z_t|z_{t-1}, x) dz_0, \dots, dz_{T-1}$$

- ▶ Variational lower bound

$$L = \mathbb{E}_{q(z_T|x)} \log \frac{p(x, z_T)}{q(z_T|x)} \leq \log p(x)$$

- ▶ The stochastic Markov chain, therefore, can be viewed as a nonparametric variational approximation.
- ▶ Can we combine MCMC and VI to get the best of both worlds?

- ▶ Use auxiliary random variables $y = (z_0, \dots, z_{T-1})$ to construct a tractable lower bound

$$L_{\text{aux}} = \mathbb{E}_{q(y, z_T|x)} \log \frac{p(x, z_T)r(y|z_T, x)}{q(y, z_T|x)} \leq \log p(x)$$

- ▶ $r(y|z_T, x)$ is an arbitrary **auxiliary distribution**, e.g.

$$r(y|z_T, x) = \prod_{t=1}^T r_t(z_{t-1}|z_t, x)$$

- ▶ This is a looser lower bound

$$\begin{aligned} L_{\text{aux}} &= \mathbb{E}_{q(y, z_T|x)} (\log p(x, z_T) + \log r(y|z_T, x) - \log q(y, z_T|x)) \\ &= L - \mathbb{E}_{q(z_T|x)} (D_{KL}(q(y|z_T, x) \| r(y|z_T, x))) \\ &\leq L \leq \log p(x) \end{aligned}$$

- ▶ Suppose z_0, z_1, \dots, z_T is a sampled trajectory

$$z_0 \sim q(z_0|x)$$

$$z_t \sim q_t(z_t|z_{t-1}, x), \quad t = 1, \dots, T$$

- ▶ Unbiased stochastic estimate of L_{aux}

$$\begin{aligned}\hat{L}_{\text{aux}} &= \log p(x, z_T) - \log q(z_0|x) + \sum_{t=1}^T \left(\log \frac{r_t(z_{t-1}|z_t, x)}{q_t(z_t|z_{t-1}, x)} \right) \\ &= \log p(x, z_0) - \log q(z_0|x) + \sum_{t=1}^T \log \alpha_t\end{aligned}$$

where

$$\alpha_t = \frac{p(x, z_t) r_t(z_{t-1}|z_t, x)}{p(x, z_{t-1}) q_t(z_t|z_{t-1}, x)}$$

- ▶ Using the detailed balance condition

$$\alpha_t = \frac{p(x, z_t) r_t(z_{t-1} | z_t, x)}{p(x, z_{t-1}) q_t(z_t | z_{t-1}, x)} = \frac{r_t(z_{t-1} | z_t, x)}{q_t(z_{t-1} | z_t, x)}$$

- ▶ Therefore,

$$L_{\text{aux}} = \mathbb{E}_{q(z_0|x)} \log \frac{p(x, z_0)}{q(z_0|x)} + \sum_{t=1}^T \mathbb{E}_{q(y, z_T|x)} \log \frac{r_t(z_{t-1} | z_t, x)}{q_t(z_{t-1} | z_t, x)}$$

- ▶ For optimal $r_t(z_{t-1} | z_t, x) = q(z_{t-1} | z_t, x)$

$$\mathbb{E}_q \log \frac{r_t(z_{t-1} | z_t, x)}{q_t(z_{t-1} | z_t, x)} = \mathbb{E}_q \log \frac{q(z_{t-1} | z_t, x)}{q_t(z_{t-1} | z_t, x)} \geq 0$$

- ▶ MCMC iterations always improve approximation unless already perfect! In practice, we need

$$r_t(z_{t-1} | z_t, x) \approx q(z_{t-1} | z_t, x)$$

- ▶ Specify a parameterized Markov chain

$$q_{\theta}(z) = q_{\theta}(z_0|x) \prod_{t=1}^T q_{\theta}(z_t|z_{t-1}, x)$$

- ▶ Specify a parameterized auxiliary distribution $r_{\theta}(y|z_T, x)$
- ▶ Sample MCMC trajectories for the variational lower bound

$$\hat{L}(\theta) = \log p(x, z_T) - \log q(z_0|x) + \sum_{t=1}^T \left(\log \frac{r_t(z_{t-1}|z_t, x)}{q_t(z_t|z_{t-1}, x)} \right)$$

- ▶ Run SGD using $\nabla_{\theta}\hat{L}(\theta)$ (reparameterization trick)

- ▶ A bivariate Gaussian target distribution

$$p(z^1, z^2) \propto \exp\left(-\frac{1}{2\tau_1^2}(z^1 - z^2)^2 - \frac{1}{2\tau_2^2}(z^1 + z^2)^2\right)$$

- ▶ Gibbs sampling

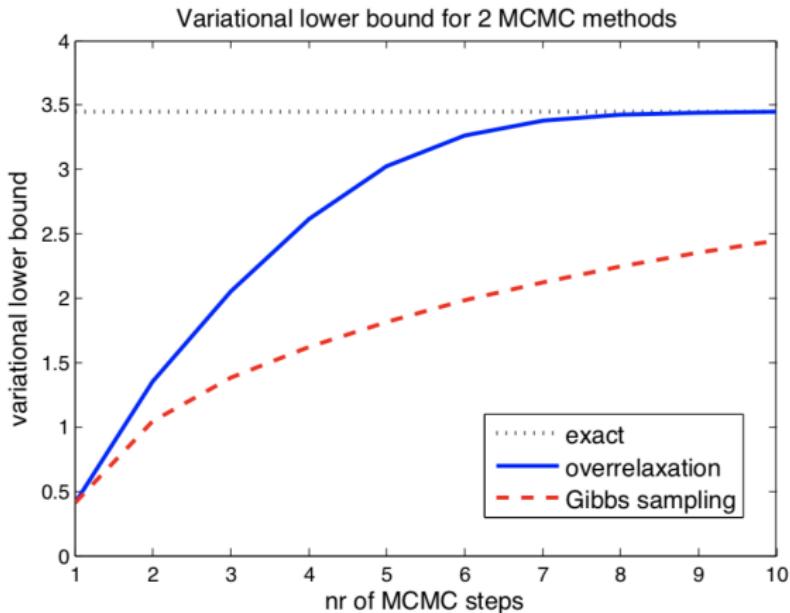
$$q(z_t^i | z_{t-1}) = p(z^i | z^{-i}) = \mathcal{N}(\mu_i, \sigma_i^2)$$

- ▶ Over-relaxation (Adler, 1981)

$$q(z_t^i | z_{t-1}) = \mathcal{N}(\mu_i + \alpha(z_{t-1}^i - \mu_i), \sigma_i^2(1 - \alpha^2))$$

- ▶ Gaussian reverse model $r_t(z_{t-1}|z_t)$, linear dependence on z_t .
Find the best α via variational lower bound maximization.

Gibbs sampling versus over-relaxation for a bivariate Gaussian



The improved mixing of over-relaxation results in an improved variational lower bound.

- We can use Hamiltonian dynamics for more efficient transition distributions

$$v'_t \sim q(v'_t | z_{t-1}, x), \quad (v_t, z_t) = \Phi(v'_t, z_{t-1})$$

where $\Phi : \mathbb{R}^{2n} \mapsto \mathbb{R}^{2n}$ is the Hamiltonian flow.

- Φ is deterministic, invertible and volume preserving

$$q(v_t, z_t | z_{t-1}, x) = q(v'_t | z_{t-1}, x), \quad r(v'_t, z_{t-1} | z_t, x) = r(v_t | z_t, x)$$

- Note that we would use *leapfrog* integrator to discretize the Hamiltonian flow. However, the resulting map $\hat{\Phi}$ is also invertible and volume preserving, and the above equations still hold.

► HMC trajectory

$$z_0 \sim q(z_0|x)$$

$$v'_t \sim q_t(v'_t|z_{t-1}, x), \quad v_t, z_t = \hat{\Phi}(v'_t, z_{t-1}), \quad t = 1, \dots, T$$

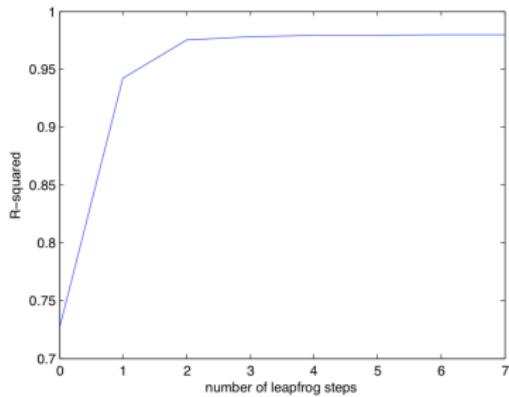
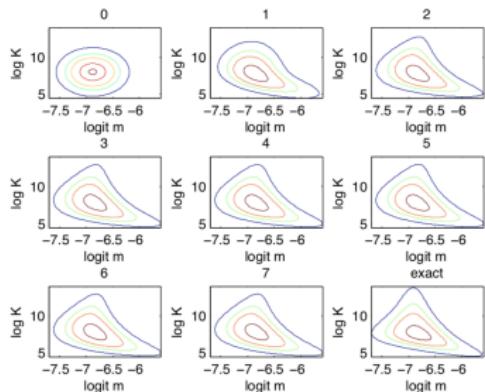
► Lower bound estimate

$$\hat{L}(\theta) = \log p(x, z_0) - \log q(z_0|x) + \sum_{t=1}^T \log \frac{p(x, z_t) r_t(v_t|z_t, x)}{p(x, z_{t-1}) q_t(v'_t|x, z_{t-1})}$$

► Stochastic optimization using $\nabla_\theta \hat{L}(\theta)$

- No rejection step, to keep everything differentiable.
- θ includes all parameters in q and r , and may include some HMC hyperparameters (stepsize and mass matrix) as well.
- Differentiate through the leapfrog integrator.

A simple 2-dimensional beta-binomial model for overdispersion.
One step of Hamiltonian dynamics with varying number of leapfrog steps.



Variational autoencoder for binarized MNIST, Gaussian prior $p(z) = \mathcal{N}(0, I)$, MLP conditional likelihood $p_\theta(x|z)$

Model	$-L$	$-\log p(x)$
<i>Results with $q(z_0 x) = \mathcal{N}(\mu, \sigma^2 \mathbf{I})$:</i>		
5 leapfrog steps	90.86	87.16
10 leapfrog steps	87.60	85.56
<i>With $q(z_0 x) = \text{inference network}$:</i>		
No leapfrog steps	94.18	88.95
1 leapfrog step	91.70	88.08
4 leapfrog steps	89.82	86.40
8 leapfrog steps	88.30	85.51

- MCMC makes bound tighter, give better marginal likelihood.
- MCMC also works with simple initialization.

Variational autoencoder for binarized MNIST, Gaussian prior $p(z) = \mathcal{N}(0, I)$, MLP conditional likelihood $p_\theta(x|z)$

Model	$-L$	$-\log p(x)$
<i>Results with $q(z_0 x) = \mathcal{N}(\mu, \sigma^2 \mathbf{I})$:</i>		
5 leapfrog steps	90.86	87.16
10 leapfrog steps	87.60	85.56
<i>With $q(z_0 x) = \text{inference network}$:</i>		
No leapfrog steps	94.18	88.95
1 leapfrog step	91.70	88.08
4 leapfrog steps	89.82	86.40
8 leapfrog steps	88.30	85.51

- MCMC makes bound tighter, give better marginal likelihood.
- MCMC also works with simple initialization.

Variational autoencoder for binarized MNIST, Gaussian prior $p(z) = \mathcal{N}(0, I)$, MLP conditional likelihood $p_\theta(x|z)$

Model	$-L$	$-\log p(x)$
<i>Results with $q(z_0 x) = \mathcal{N}(\mu, \sigma^2 \mathbf{I})$:</i>		
5 leapfrog steps	90.86	87.16
10 leapfrog steps	87.60	85.56
<i>With $q(z_0 x) = \text{inference network}$:</i>		
No leapfrog steps	94.18	88.95
1 leapfrog step	91.70	88.08
4 leapfrog steps	89.82	86.40
8 leapfrog steps	88.30	85.51

- MCMC makes bound tighter, give better marginal likelihood.
- MCMC also works with simple initialization.

- ▶ MCMC improves variational approximation
 - ▶ MCMC kernels automatically adapt to target $p(z|x)$.
 - ▶ More flexible approximations in addition to standard exponential family distributions.
 - ▶ More MCMC steps \Rightarrow slower iterations, but few iterations needed for convergence.
- ▶ Optimizing variational bound improves MCMC
 - ▶ Automatic tuning, convergence assessment, independent sampling, no rejections.
 - ▶ Learning MCMC transitions $q_t(z_t|z_{t-1}, x)$.
 - ▶ Optimize initialization $q(z_0|x)$.
- ▶ Many possibilities left to explore.

- ▶ D. J. Rezende and S. Mohamed. Variational inference with normalizing flows. Proceedings of the 32nd International Conference on Machine Learning, pages 1530–1538, 2015.
- ▶ L. Dinh, D. Krueger, and Y. Bengio. NICE: Non-linear Independent Components Estimation. arXiv:1410.8516, 2014.
- ▶ L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using Real NVP. Proceedings of the 5th International Conference on Learning Representations, 2017.
- ▶ Stephen L Adler. Over-relaxation method for the monte carlo evaluation of the partition function for multiquadratic actions. Physical Review D, 23(12):2901, 1981.

- ▶ D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling. Improved variational inference with Inverse Autoregressive Flow. *Advances in Neural Information Processing Systems* 29, pages 4743–4751, 2016.
- ▶ Papamakarios, G., Murray, I., and Pavlakou, T. (2017). Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pages 2335–2344.
- ▶ T. Salimans,D. P. Kingma, and M. Welling. Markov chain monte carlo and variational inference: Bridging the gap. In *ICML*, 2015.