
Neural Sampling in Hierarchical Exponential-family Energy-based Models: Supplementary Information

Xingsi Dong¹

dxs19980605@pku.edu.cn

Si Wu^{1,2}

siwu@pku.edu.cn

1. PKU-Tsinghua Center for Life Sciences, Academy for Advanced Interdisciplinary Studies.
School of Psychology and Cognitive Sciences, Peking University.
2. IDG/McGovern Institute for Brain Research. Center of Quantitative Biology, Peking University.

1 Proof of equations in the main text

1.1 The derivative of Kullback–Leibler divergence

Derivation of Eq.(2) in the main text,

$$\nabla_{\theta} D_{\text{KL}} [p_{\text{true}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})], \quad (1a)$$

$$= -\mathbb{E}_{\mathbf{x} \sim p_{\text{true}}(\mathbf{x})} \frac{1}{p_{\theta}(\mathbf{x})} \nabla_{\theta} p_{\theta}(\mathbf{x}), \quad (1b)$$

$$= -\mathbb{E}_{\mathbf{x} \sim p_{\text{true}}(\mathbf{x})} \frac{1}{p_{\theta}(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z}|\mathbf{x})} \left[\frac{\nabla_{\theta} p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right], \quad (1c)$$

$$= -\mathbb{E}_{\mathbf{x} \sim p_{\text{true}}(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z}|\mathbf{x})} \left[\frac{\nabla_{\theta} p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{x}, \mathbf{z})} \right], \quad (1d)$$

$$= -\mathbb{E}_{\mathbf{x} \sim p_{\text{true}}(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \ln p_{\theta}(\mathbf{x}, \mathbf{z})]. \quad (1e)$$

1.2 The derivative of log-partition function

Considering that the sum of the probability equals to 1, we can obtain,

$$1 = \int p(\mathbf{x}_l | \mathbf{x}_{l+1}) d\mathbf{x}_l, \quad (2a)$$

$$= \int \exp [\boldsymbol{\eta}_l^T \phi(\mathbf{x}_l) + g(\mathbf{x}_l) - A(\boldsymbol{\eta}_l)] d\mathbf{x}_l, \quad (2b)$$

$$= \frac{\int \exp [\boldsymbol{\eta}_l^T \phi(\mathbf{x}_l) + g(\mathbf{x}_l)] d\mathbf{x}_l}{\exp [A(\boldsymbol{\eta}_l)]}. \quad (2c)$$

Then, we take the derivative of $\exp[A(\boldsymbol{\eta}_l)]$ with respect to $\boldsymbol{\eta}_l$, obtaining,

$$A'(\boldsymbol{\eta}_l) = \frac{(\exp[A(\boldsymbol{\eta}_l)])'}{\exp[A(\boldsymbol{\eta}_l)]}, \quad (3a)$$

$$= \frac{(\int \exp[\boldsymbol{\eta}_l^T \phi(\mathbf{x}_l) + g(\mathbf{x}_l)] d\mathbf{x}_l)'}{\exp[A(\boldsymbol{\eta}_l)]}, \quad (3b)$$

$$= \frac{\int \phi(\mathbf{x}_l) \exp[\boldsymbol{\eta}_l^T \phi(\mathbf{x}_l) + g(\mathbf{x}_l)] d\mathbf{x}_l}{\exp[A(\boldsymbol{\eta}_l)]}, \quad (3c)$$

$$= \int \phi(\mathbf{x}_l) \exp[\boldsymbol{\eta}_l^T \phi(\mathbf{x}_l) + g(\mathbf{x}_l) - A(\boldsymbol{\eta}_l)] d\mathbf{x}_l, \quad (3d)$$

$$= E_{\mathbf{x}_l \sim p_\theta(\mathbf{x}_l | \mathbf{x}_{l+})} [\phi(\mathbf{x}_l)]. \quad (3e)$$

1.3 Second-order Langevin dynamic

In this section, we provide a proof that the stationary distribution $\tilde{p}(\mathbf{z})$ of the second-order Langevin dynamic, as described in the main text, is equivalent to the posterior distribution $p_\theta(\mathbf{z} | \mathbf{x})$.

$$\tau_z \frac{d\mathbf{z}}{dt} = \nabla_{\mathbf{z}} \ln p_\theta(\mathbf{z} | \mathbf{x}) - \mathbf{v} + \sqrt{2\tau_z} \xi, \quad (4)$$

$$\tau_v \frac{d\mathbf{v}}{dt} = -\frac{m\mathbf{v}}{2} + m\sqrt{2\tau_v} \xi, \quad (5)$$

where ξ is Gaussian white noise satisfying $\langle \xi \xi^T \rangle = \mathbf{I} \delta(t - t')$. Following the anti-symmetric matrix decomposition strategy [1], the above dynamic can be reorganised to,

$$\frac{d}{dt} \begin{pmatrix} \mathbf{z} \\ \mathbf{v} \end{pmatrix} = (D + Q) \begin{pmatrix} \nabla_{\mathbf{z}} \ln p_\theta(\mathbf{z} | \mathbf{x}) \\ -\tau_v \mathbf{v} / (4m\tau_z) \end{pmatrix} + \sqrt{2D} \xi, \quad (6)$$

where ξ is Gaussian white noise satisfying $\langle \xi \xi^T \rangle = \mathbf{I} \delta(t - t')$ and

$$D = \begin{pmatrix} I/\tau_z & 2mI\tau_v \\ 2m/\tau_v & 4m^2\tau_z/\tau_v^2 \end{pmatrix}, \quad Q = \begin{pmatrix} 0 & 2mI\tau_v \\ -2m/\tau_v & 0 \end{pmatrix}. \quad (7)$$

The Fokker-Planck equation of the above dynamics is written as,

$$\frac{\partial}{\partial t} p(\mathbf{z}, \mathbf{v}, t) = \nabla^\top \left[(D + Q) \begin{pmatrix} \nabla_{\mathbf{z}} \ln p_\theta(\mathbf{z} | \mathbf{x}) \\ -\tau_v \mathbf{v} / (4m\tau_z) \end{pmatrix} p(\mathbf{z}, \mathbf{v}, t) \right] + \nabla^\top [D \nabla p(\mathbf{z}, \mathbf{v}, t)]. \quad (8)$$

By utilizing the property,

$$\nabla^\top [Q \nabla p(\mathbf{z}, \mathbf{v}, t)] = 0, \quad (9)$$

Eq.(8) can be rewritten as,

$$\frac{\partial}{\partial t} p(\mathbf{z}, \mathbf{v}, t) = \nabla^\top \left\{ (D + Q) \left[\begin{pmatrix} \nabla_{\mathbf{z}} \ln p_\theta(\mathbf{z} | \mathbf{x}) \\ -\tau_v \mathbf{v} / (4m\tau_z) \end{pmatrix} p(\mathbf{z}, \mathbf{v}, t) + \nabla p(\mathbf{z}, \mathbf{v}, t) \right] \right\}. \quad (10)$$

By letting the left-hand side of the above equation to be zero, we get that the stationary distribution $\tilde{p}(\mathbf{z}, \mathbf{v})$ satisfy,

$$\begin{pmatrix} \nabla_{\mathbf{z}} \ln p_\theta(\mathbf{z} | \mathbf{x}) \\ -\tau_v \mathbf{v} / (4m\tau_z) \end{pmatrix} p(\mathbf{z}, \mathbf{v}) + \nabla p(\mathbf{z}, \mathbf{v}) = 0, \quad (11)$$

which gives the solution,

$$\tilde{p}(\mathbf{z}, \mathbf{v}) = p_\theta(\mathbf{z} | \mathbf{x}) \mathcal{N}(\mathbf{v}; 0, 4m\tau_z/\tau_v). \quad (12)$$

The marginal stationary distribution of \mathbf{z} is then calculated to be,

$$\tilde{p}(\mathbf{z}) = \int \tilde{p}(\mathbf{z}, \mathbf{v}) d\mathbf{v} = p_\theta(\mathbf{z} | \mathbf{x}). \quad (13)$$

2 Analysis of Hessian matrix

In this section, we provide a theoretical lower bound for $\lambda_1(H)$ and an upper bound for $\det(H)$ based on the assumption that $\phi(\mathbf{x}_l) = \mathbf{x}_l$ and $g(\mathbf{x}_l) = -\mathbf{x}_l^T \mathbf{x}_l$. Then we demonstrate that HEE-L can only converge to unimodal distributions.

We define the total energy $F_t = -\ln p_\theta(\mathbf{x}_{0:L})$. Considering

$$\frac{\partial^2}{\partial \mathbf{x}_l \partial \mathbf{x}_{l+k}} F_t = \begin{cases} I + \theta_{l-1}^T \theta_{l-1} & , k = 0 \\ -\theta_l & , k = 1 \\ -\theta_{l-1} & , k = -1 \\ 0 & , |k| > 1 \end{cases} \quad (14)$$

therefore, the Hessian matrix $H_t = \Delta_{\mathbf{x}_{0:L}} F_t$ is calculated as,

$$H_t = (I + D)^T (I + D) \quad (15)$$

where

$$D = \begin{pmatrix} 0 & -\theta_0 & & & \\ & 0 & -\theta_1 & & \\ & & \ddots & \ddots & \\ & & & 0 & -\theta_{L-1} \\ & & & & 0 \end{pmatrix} \quad (16)$$

Thus, the lower bound of singular values of $I + D$ is given by (see Theorem 8.13 in [2]),

$$\sigma_1(I + D) \geq \sigma_1(I) - \sigma_{\max}(D) \quad (17)$$

where $\sigma_1(\cdot)$ and $\sigma_{\max}(\cdot)$ denote the smallest and greatest singular value, respectively. Then we assume that θ_l is matrix with sub-gaussian entries, and $n_1 = \dots = n_L = N/L$ then we have,

$$\begin{aligned} \mathbb{E}[\sigma_{\max}(D)] &= \mathbb{E}[\max_{l \leq L} \sigma_{\max}(\theta_l)], \\ &\leq K \sqrt{\ln L} \mathbb{E}[\sigma_{\max}(\theta_L)], \quad \text{Exercise 2.5.10 in [3]} \\ &\leq K \sqrt{\frac{N \ln L}{L}}, \quad \text{Theorem 4.4.5 in [3]} \end{aligned} \quad (18)$$

where K is a small constant. Thus, the lower bound of $\lambda_1(H_t)$ is calculated as

$$\begin{aligned} \lambda_1(H_t) &= \sigma_1^2(I + D) \quad \text{Eq.(15)} \\ &\geq \left(1 - K \sqrt{\frac{N \ln L}{L}}\right)^2 \quad \text{Eq.(18)} \end{aligned} \quad (19)$$

Thus, the lower bound of $\lambda_1(H)$ is calculated as (Theorem 4.3.15 in [4]),

$$\lambda_1(H) \geq \lambda_1(H_t). \quad (20)$$

And the upper bound of $\det(H)$ is calculated as,

$$\begin{aligned} \det(H) &= \prod_{i=1}^N \lambda_i(H) \\ &\leq \prod_{i=1}^N \lambda_{i+n_0}(H_t) \quad \text{Theorem 4.3.15 in [4]} \\ &= 1 / \prod_{i=1}^{n_0} \lambda_i(H_t) \quad \text{Utilizing } \det(H_t) = 1 \\ &\leq 1 / [\lambda_1(H_t)]^{n_0} \end{aligned} \quad (21)$$

The above analysis reveals that $\lambda_1(H_t)$ increase with the decrease of L . Additionally, as the value of $\lambda_1(H_t)$ increases, the lower bound of $\lambda_1(H)$ and the upper bound of $\det(H)$ will decrease and increase, respectively.

Then, we prove that the Hessian matrix $H_0 = -\Delta_{\mathbf{x}_0} \ln p_\theta(\mathbf{x}_0)$ is positive-definite, indicating that HEE-L is capable of approximating only unimodal distributions.

$$H_0 = -\Delta_{\mathbf{x}_0} \ln p_\theta(\mathbf{x}_0), \quad (22a)$$

$$= -\Delta_{\mathbf{x}_0} \ln \int p_\theta(\mathbf{x}_{0:L}) d\mathbf{x}_{1:L}, \quad (22b)$$

$$= \frac{-1}{[p_\theta(\mathbf{x}_0)]^2} \int d\mathbf{x}'_{1:L} \int d\mathbf{x}_{1:L} \left[p_\theta(\mathbf{x}'_{0:L}) \Delta_{\mathbf{x}_0} p_\theta(\mathbf{x}_{0:L}) - \nabla_{\mathbf{x}_0} p_\theta(\mathbf{x}_{0:L}) \nabla_{\mathbf{x}'_0}^T p_\theta(\mathbf{x}'_{0:L}) \right] \quad (22c)$$

$$= \frac{-2}{[p_\theta(\mathbf{x}_0)]^2} \iint d\mathbf{x}_{1:L} d\mathbf{x}'_{1:L} [P(\mathbf{x}_{0:L}, \mathbf{x}'_{0:L}) + P(\mathbf{x}'_{0:L}, \mathbf{x}_{0:L})], \quad (22d)$$

$$\succ \frac{-2}{[p_\theta(\mathbf{x}_0)]^2} \iint d\mathbf{x}_{1:L} d\mathbf{x}'_{1:L} [P(\mathbf{x}_{0:L}, \mathbf{x}_{0:L}) + P(\mathbf{x}'_{0:L}, \mathbf{x}'_{0:L})], \quad (22e)$$

$$= \frac{-1}{[p_\theta(\mathbf{x}_0)]^2} \iint d\mathbf{x}_{1:L} d\mathbf{x}'_{1:L} [p_\theta(\mathbf{x}_{0:L})]^2 \Delta_{\mathbf{x}_0} \ln p_\theta(\mathbf{x}_{0:L}), \quad (22f)$$

$$\succ 0 \quad (22g)$$

where $P(\mathbf{x}_{0:L}, \mathbf{x}'_{0:L}) = p_\theta(\mathbf{x}'_{0:L}) \Delta_{\mathbf{x}_0} p_\theta(\mathbf{x}_{0:L}) - \nabla_{\mathbf{x}_0} p_\theta(\mathbf{x}_{0:L}) \nabla_{\mathbf{x}'_0}^T p_\theta(\mathbf{x}'_{0:L})$. In Eq.(22e), we use the property

$$\Delta_{\mathbf{x}_0} p_\theta(\mathbf{x}_{0:L}) = \Delta_{\mathbf{x}'_0} p_\theta(\mathbf{x}'_{0:L}), \quad (23a)$$

$$\nabla_{\mathbf{x}_0} p_\theta(\mathbf{x}_{0:L}) \nabla_{\mathbf{x}_0}^T p_\theta(\mathbf{x}_{0:L}) + \nabla_{\mathbf{x}'_0} p_\theta(\mathbf{x}'_{0:L}) \nabla_{\mathbf{x}'_0}^T p_\theta(\mathbf{x}'_{0:L}) \succ 2 \nabla_{\mathbf{x}_0} p_\theta(\mathbf{x}_{0:L}) \nabla_{\mathbf{x}'_0}^T p_\theta(\mathbf{x}'_{0:L}) \quad (23b)$$

In Eq.(22g), we utilize the property that $-\Delta_{\mathbf{x}_0} \ln p_\theta(\mathbf{x}_{0:L})$ is a principal submatrix of H_t , which is also a positive-definite matrix.

3 Experimental details

Table 1 lists experimental settings for *2D synthetic datasets*, *FashionMNIST* and *CIFAR10*.

Settings	<i>2D synthetic datasets</i>	<i>FashionMNIST</i>	<i>CIFAR10</i>
$f(\mathbf{x}_l)$	\mathbf{x}_l	\mathbf{x}_l	\mathbf{x}_l
$g(\mathbf{x}_l)$	$-\mathbf{x}_l^T \mathbf{x}_l$	$-\mathbf{x}_l^T \mathbf{x}_l$	$-\mathbf{x}_l^T W_l \mathbf{x}_l - \text{Tr}(W_l^T W_l)$
L	5	5	10
n_l	$\text{linspace}(10, 10, L)$	$\text{linspace}(1000, 100, L)$	$\text{linspace}(3000, 100, L)$
τ_z	$\text{linspace}(1, 10, L)$	$\text{linspace}(1, 10, L)$	$\text{linspace}(10, 10, L)$
τ_x	1	1	10
τ_v	$0.5 * \tau_z$	$0.5 * \tau_z$	$0.5 * \tau_z$
m	0.5	0.5	0.5
τ_u	0.1	0.1	0.1
τ_θ	10	10	100
sparse connection	False	False	True (20%)
parameters	400	400K	4M

Table 1: Experimental settings

We conducted all experiments on a GPU (RTX A6000). The inference, learning, and generation dynamics were simulated using the Euler expansion with a step size of $dt = 0.01\tau_x$. In the inference and learning phase, each data point was presented to the network for $300\tau_x$, with each batch consisting of 30K steps. The *FashionMNIST* model was trained with a batch size of 512, taking approximately 0.5 hours to complete one epoch. The *CIFAR10* model was trained with a batch size of 128, taking approximately 2 hours to complete one epoch. And the parameters W_l and \mathbf{b}_l adopts the gradient method described by Eq.(4) in the main text.

For joint generation, we randomly initialized the neuronal activities and performed the joint generation dynamic for $300\tau_x$. The average value of neuronal activity \mathbf{x}_0 in the last $100\tau_x$ was used as the

generated output. For marginal generation, we also randomly initialized the neuronal activities and performed the marginal generation dynamic for $100\tau_x$. The average value of neuronal activity \mathbf{x}_0 in the last $50\tau_x$ was used as the generated output.

References

- [1] Yi-An Ma, Tianqi Chen, and Emily Fox. A complete recipe for stochastic gradient mcmc. *Advances in neural information processing systems*, 28, 2015.
- [2] Fuzhen Zhang. *Matrix theory: basic results and techniques*. Springer, 2011.
- [3] Omiros Papaspiliopoulos. High-dimensional probability: An introduction with applications in data science, 2020.
- [4] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.