A One-dimensional bimodal target



(a) Target (blue curve) & base (green curve) density functions. (b) Joint energy (contour plot) & example trajectory (green curve).



(c) Joint density (contour plot) and CT HMC samples (circles). (d) Histograms from HMC (top) and CT HMC (bottom) samples.

Figure 4: Visualisations of *continuous tempering* (CT) in a bimodal univariate target density. (a) A two-component Gaussian mixture target density (blue curve) and Gaussian base density (green curve) with mean and variance matched to the target. (b) The extended potential energy on the target state x and temperature control variable u (contour plot - dark colours indicate low energy) and an example simulated Hamiltonian trajectory in the joint space (green curve). The temperature control variable bridges the base and target densities lowering energy barriers in the target space. (c) Joint density on the target state x and inverse temperature β (17) (contour plot, dark colours indicate high density) and samples from a CT HMC chain run in the joint space (circles, size of each circle is proportional to $p_{\beta|x}(1 \mid x)$ and so larger symbols indicate a greater weighting in estimates of expectations with respect to the target (22)). (d) Example target state sample histograms from running standard HMC in the original target density (top) and running HMC in the extended joint space (bottom).

We give here an illustrative example of the gains of the proposed approach over standard HMC in target densities with isolated modes. We use a one-dimensional Gaussian mixture density with two separated Gaussian components as the target density, shown by the blue curve in Figure 4a. Although performance in this toy univariate model is not necessarily reflective of that in more realistic higher-dimensional models, it has the advantage of allowing the joint density on x and $\beta = \beta(u)$ to be directly visualised.

For the base density $\exp[-\psi(x)]$ we use a univariate Gaussian with mean and variance matched to those of the target density (corresponding to the Gaussian density minimising the KL divergence from the target to base distribution), shown by the green curve in Figure 4a. We also set $\log \zeta = \log Z$ and so the performance here represents a 'best-case' scenario for the continuous tempering approach.

The resulting potential energy $(-\log p_{x,u})$ on the extended (x, u) space is shown in Figure 4b. For positive temperature control values (and so inverse temperature values close to 1), the energy surface tends increasingly to the double-well potential corresponding to the target distribution, with a high energy barrier between the two modes. For negative temperature control values the energy surface tends towards the single quadratic well corresponding to the Gaussian base density. The resulting joint energy surface allows for paths between the values of the target state x corresponding to the two modes in the original target space, allowing simulated Hamiltonian trajectories such as that shown in green to more easily explore the target state space.

Samples from a HMC chain on the extended joint space are shown in Figure 4c, with the joint density on (x, β) (17) shown in the background as a contoured heat map. It can be seen that the Hamiltonian dynamic is able to explore the joint space well with good coverage of all of the high density regions. The size of the points in 4c is proportional to $w_1(x) = p_{\beta|x}(1 \mid x)$ and so reflects the importance weights of the samples in the estimator for expectations with respect to the target in (22). Importantly even points for which β is close to zero can contribute significantly to the expectations if the corresponding x value is probable under the target: this is in contrast to the extended Hamiltonian approach of [19] where only a subset of points corresponding to $\beta = 1$ are used to compute expectations.

The final panel, Figure 4d shows empirical histograms on the target variable x estimated from samples of a chain on the extended space (joint continuous tempering, bottom) and standard HMC on the original target space (top). As can be seen the standard HMC approach gets stuck in one mode thus does not assign any mass to the other mode in the histogram, unlike the tempered chain which identifies both modes and accurately estimates their relative masses.

B Bounding the inverse temperature marginal density

We have a joint density on (\mathbf{x}, β)

$$p_{\mathbf{x},\beta}(\mathbf{x},\beta) = \frac{1}{C} \exp(-\beta \phi(\mathbf{x}) - \beta \log \zeta - (1-\beta)\psi(\mathbf{x})).$$
(27)

The resulting marginal density on β is

$$p_{\beta}(\beta) = \int_{\mathcal{X}} p_{\mathbf{x},\beta}(\mathbf{x},\beta) \, \mathrm{d}\mathbf{x} = \frac{1}{C\zeta^{\beta}} \int_{\mathcal{X}} \exp(-\beta\phi(\mathbf{x}) - (1-\beta)\psi(\mathbf{x})) \, \mathrm{d}\mathbf{x}.$$
 (28)

To derive an upper-bound on $p_{\beta}(\beta)$ we use Hölder's inequality

$$\int_{\mathcal{X}} g(\mathbf{x}) h(\mathbf{x}) \, \mathrm{d}\mathbf{x} \le \left(\int_{\mathcal{X}} |g(\mathbf{x})|^{\frac{1}{a}} \, \mathrm{d}\mathbf{x} \right)^{a} \left(\int_{\mathcal{X}} |h(\mathbf{x})|^{\frac{1}{1-a}} \, \mathrm{d}\mathbf{x} \right)^{1-a}$$
(29)

where $a \in [0, 1]$ and g and h are measurable functions. We also use the definitions

$$\int_{\mathcal{X}} \exp(-\phi(\mathbf{x})) \, \mathrm{d}\mathbf{x} = Z \quad \text{and} \quad \int_{\mathcal{X}} \exp(-\psi(\mathbf{x})) \, \mathrm{d}\mathbf{x} = 1.$$
(30)

From (28) we have that

$$\mathsf{p}_{\beta}(\beta) = \frac{1}{C\zeta^{\beta}} \int_{\mathcal{X}} \left(\exp(-\phi(\mathbf{x}))^{\beta} \right) \left(\exp(-\psi(\mathbf{x}))^{1-\beta} \right) \mathrm{d}\mathbf{x}.$$
(31)

Applying Hölder's inequality (29) with $g(\mathbf{x}) = \exp(-\phi(\mathbf{x}))^{\beta}$, $h(\mathbf{x}) = \exp(-\psi(\mathbf{x}))^{1-\beta}$ and $a = \beta$

$$p_{\beta}(\beta) \leq \frac{1}{C\zeta^{\beta}} \left(\int_{\mathcal{X}} \left| \exp(-\phi(\mathbf{x}))^{\beta} \right|^{\frac{1}{\beta}} d\mathbf{x} \right)^{\beta} \left(\int_{\mathcal{X}} \left| \exp(-\psi(\mathbf{x}))^{1-\beta} \right|^{\frac{1}{1-\beta}} d\mathbf{x} \right)^{1-\beta}$$
(32)

$$= \frac{1}{C\zeta^{\beta}} \left(\int_{\mathcal{X}} \exp(-\phi(\mathbf{x})) \, \mathrm{d}\mathbf{x} \right)^{\rho} \left(\int_{\mathcal{X}} \exp(-\psi(\mathbf{x})) \, \mathrm{d}\mathbf{x} \right)^{1-\rho}.$$
 (33)

Substituting the definitions in (30) gives

$$\mathsf{p}_{\beta}(\beta) \le \frac{1}{C} \left(\frac{Z}{\zeta}\right)^{\beta}.$$
(34)

To derive a lower-bound on $p_{\beta}(\beta)$, we use Jensen's inequality

$$\varphi\left(\int_{\mathcal{X}} g(\mathbf{x})q(\mathbf{x})\,\mathrm{d}\mathbf{x}\right) \ge \int_{\mathcal{X}} \varphi(g(\mathbf{x}))q(\mathbf{x})\,\mathrm{d}\mathbf{x},\tag{35}$$

for a concave function φ , normalised density $q : \int_{\mathcal{X}} q(\mathbf{x}) d\mathbf{x} = 1$ and measurable g. The logarithm of (28) gives

$$\log p_{\beta}(\beta) + \beta \log \zeta + \log C = \log \left(\int_{\mathcal{X}} \exp(-\beta(\phi(\mathbf{x}) - \psi(\mathbf{x}))) \exp(-\psi(\mathbf{x})) \, \mathrm{d}\mathbf{x} \right).$$
(36)

Applying Jensen's inequality (35) with $\varphi = \log_{10} q = \exp(-\psi)$ and $g = \exp(-\beta(\phi - \psi))$

$$\log p_{\beta}(\beta) + \log C + \beta \log \zeta \ge \beta \int_{\mathcal{X}} (\psi(\mathbf{x}) - \phi(\mathbf{x})) \exp(-\psi(\mathbf{x})) \, \mathrm{d}\mathbf{x}$$
(37)

$$= \beta \int_{\mathcal{X}} (\log Z - \log Z - \log \exp(-\psi(\mathbf{x}) + \phi(\mathbf{x}))) \exp(-\psi(\mathbf{x})) d\mathbf{x}$$
(38)

$$= \beta \log Z - \beta \int_{\mathcal{X}} \exp(-\psi(\mathbf{x})) \log\left(\frac{\exp(-\psi(\mathbf{x}))}{\exp(-\phi(\mathbf{x}))/Z}\right) d\mathbf{x}.$$
 (39)

Recognising the integral in the last line as the *Kullback–Leibler* (KL) divergence $d^{b \to t}$ from the base density $\exp(-\psi(\mathbf{x}))$ to the target density $\exp(-\phi(\mathbf{x}))/Z$

$$d^{b \to t} = \int_{\mathcal{X}} \exp(-\psi(\mathbf{x})) \log\left(\frac{\exp(-\psi(\mathbf{x}))}{\exp(-\phi(\mathbf{x}))/Z}\right) d\mathbf{x},\tag{40}$$

and taking the exponential of both sides and rearranging we have

$$\mathsf{p}_{\beta}(\beta) \ge \frac{1}{C} \left(\frac{Z}{\zeta}\right)^{\beta} \exp\left(-\beta d^{b \to t}\right). \tag{41}$$

By instead noting (28) can be rearranged into the form

$$\log p_{\beta}(\beta) + \log C + \beta \log \zeta - \log Z = \log \left(\int_{\mathcal{X}} \exp(-(1-\beta)(\psi(\boldsymbol{x}) - \phi(\boldsymbol{x}))) \frac{1}{Z} \exp(-\phi(\boldsymbol{x})) \, \mathrm{d}\boldsymbol{x} \right), \tag{42}$$

by an equivalent series of steps we can also derive a bound using the reversed form of the KL divergence

$$d^{t \to b} = \int_{\mathcal{X}} \frac{1}{Z} \exp(-\phi(\mathbf{x})) \log\left(\frac{\exp(-\phi(\mathbf{x}))/Z}{\exp(-\psi(\mathbf{x}))}\right) d\mathbf{x}.$$
 (43)

from the target to the base distribution, giving that

$$\mathsf{p}_{\beta}(\beta) \ge \frac{1}{C} \left(\frac{Z}{\zeta}\right)^{\beta} \exp\left(-(1-\beta)d^{t \to b}\right). \tag{44}$$

C Gaussian mixture Boltzmann machine relaxations

We define a *Boltzmann machine distribution* on a signed binary state $\mathbf{s} \in \{-1, +1\}^{D_B} = S$ as

$$\mathsf{p}_{\mathbf{s}}(\mathbf{s}) = \frac{1}{Z_B} \exp\left(\frac{1}{2}\mathbf{s}^\mathsf{T} \mathbf{W} \mathbf{s} + \mathbf{s}^\mathsf{T} \mathbf{b}\right) \qquad Z_B = \sum_{\mathbf{s} \in S} \left(\exp\left(\frac{1}{2}\mathbf{s}^\mathsf{T} \mathbf{W} \mathbf{s} + \mathbf{s}^\mathsf{T} \mathbf{b}\right)\right). \tag{45}$$

We introduce an auxiliary real-valued vector random variable $\mathbf{x} \in \mathbb{R}^{D}$ with a Gaussian conditional distribution

$$p_{\mathbf{x}|\mathbf{s}}(\mathbf{x} \mid \mathbf{s}) = \frac{1}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2} \left(\mathbf{x} - \mathbf{Q}^{\mathsf{T}} \mathbf{s}\right)^{\mathsf{T}} \left(\mathbf{x} - \mathbf{Q}^{\mathsf{T}} \mathbf{s}\right)\right)$$
(46)

with $Q \ a \ D_B \times D$ matrix such that $QQ^T = W + D$ for some diagonal D which makes W + D positive semi-definite. In our experiments, based on the observation in [44] that minimising the maximum eigenvalue of W + D decreases the maximal separation between the Gaussian components in the relaxation, we set D as the solution to the semi-definite programme

$$\min_{\boldsymbol{D}} \left(\lambda_{\text{MAX}} (\boldsymbol{W} + \boldsymbol{D}) \right) : \boldsymbol{W} + \boldsymbol{D} \ge 0$$
(47)

where λ_{MAX} denotes the maximal eigenvalue. In general the optimised W + D lies on the semi-definite cone and so has rank less than D_B hence a Q can be found such that $D < D_B$. The resulting joint distribution on (\mathbf{x}, \mathbf{s}) is

$$\mathsf{p}_{\mathbf{x},\mathbf{s}}(\mathbf{x},\mathbf{s}) = \frac{1}{(2\pi)^{D/2} Z_B} \exp\left(-\frac{1}{2}\mathbf{x}^\mathsf{T}\mathbf{x} + \mathbf{s}^\mathsf{T}\mathbf{Q}\mathbf{x} - \frac{1}{2}\mathbf{s}^\mathsf{T}\mathbf{Q}\mathbf{Q}^\mathsf{T}\mathbf{s} + \frac{1}{2}\mathbf{s}^\mathsf{T}\mathbf{W}\mathbf{s} + \mathbf{s}^\mathsf{T}\mathbf{b}\right)$$
(48)

$$= \frac{1}{(2\pi)^{D/2} Z_B} \exp\left(-\frac{1}{2} \mathbf{x}^\mathsf{T} \mathbf{x} + \mathbf{s}^\mathsf{T} (\mathbf{Q} \mathbf{x} + \mathbf{b}) - \frac{1}{2} \mathbf{s}^\mathsf{T} \mathbf{D} \mathbf{s}\right)$$
(49)

$$= \frac{1}{(2\pi)^{D/2} Z_B \exp\left(\frac{1}{2} \operatorname{Tr}(\boldsymbol{D})\right)} \exp\left(-\frac{1}{2} \boldsymbol{x}^{\mathsf{T}} \boldsymbol{x}\right) \prod_{i=1}^{D_B} \left(\exp\left(s_i \left(\boldsymbol{q}_i^{\mathsf{T}} \boldsymbol{x} + \boldsymbol{b}_i\right)\right)\right),\tag{50}$$

where $\{q_i^{\mathsf{T}}\}_{i=1}^{D_B}$ are the D_B rows of Q. We can marginalise over the binary state **s** as each s_i is conditionally independent of all the others given **x** in the joint distribution. This gives the *Boltzmann machine relaxation* density on **x**

$$\mathsf{p}_{\mathbf{x}}(\mathbf{x}) = \frac{2^{D_B}}{(2\pi)^{D/2} Z_B \exp\left(\frac{1}{2} \operatorname{Tr}(\mathbf{D})\right)} \exp\left(-\frac{1}{2} \mathbf{x}^\mathsf{T} \mathbf{x}\right) \prod_{i=1}^{D_B} \left(\cosh\left(\mathbf{q}_i^\mathsf{T} \mathbf{x} + b_i\right)\right),\tag{51}$$

which is a structured Gaussian mixture density with 2^{D_B} components. If we define $p_x(x) = \frac{1}{Z} \exp(-\phi(x))$ with

$$\phi(\mathbf{x}) = \frac{1}{2}\mathbf{x}^{\mathsf{T}}\mathbf{x} - \sum_{i=1}^{D_B} \left(\log\cosh\left(\mathbf{q}_i^{\mathsf{T}}\mathbf{x} + b_i\right)\right),\tag{52}$$

then the normalisation constant Z of the relaxation density can be related to the normalising constant of the corresponding Boltzmann machine distribution by

$$\log Z = \log Z_B + \frac{1}{2} \operatorname{Tr}(D) + \frac{D}{2} \log(2\pi) - D_B \log 2.$$
(53)

It can also be shown that the first and second moments of the relaxation distribution are related to the first and second moments of the corresponding Boltzmann machine distribution by

$$\mathbb{E}[\mathbf{x}] = \int_{\mathcal{X}} \mathbf{x} \sum_{s \in S} \left(p_{\mathbf{x}|s}(\mathbf{x} \mid s) p_{s}(s) \right) d\mathbf{x} = \sum_{s \in S} \left(\int_{\mathcal{X}} \mathbf{x} \, \mathcal{N}\left(\mathbf{x}; \mathbf{Q}^{\mathsf{T}} s, \mathbf{I}\right) d\mathbf{x} \, p_{s}(s) \right) = \mathbb{E}\left[\mathbf{Q}^{\mathsf{T}} s\right] = \mathbf{Q}^{\mathsf{T}} \mathbb{E}[s], \qquad (54)$$

and
$$\mathbb{E}[\mathbf{x}\mathbf{x}^{\mathsf{T}}] = \sum_{s \in S} \left(\int_{\mathcal{X}} \mathbf{x}\mathbf{x}^{\mathsf{T}} \,\mathcal{N}(\mathbf{x}; \mathbf{Q}^{\mathsf{T}}s, \mathbf{I}) \,\mathrm{d}\mathbf{x} \,\mathrm{p}_{\mathsf{s}}(s) \right) = \mathbb{E}[\mathbf{Q}^{\mathsf{T}} \mathbf{s}\mathbf{s}\mathbf{Q} + \mathbf{I}] = \mathbf{Q}^{\mathsf{T}} \,\mathbb{E}[\mathbf{s}\mathbf{s}^{\mathsf{T}}]\mathbf{Q} + \mathbf{I}.$$
 (55)

The weight parameters W of the Boltzmann machine distributions used in the experiments in Section 7.1 were generated using an eigendecomposition based method. A uniformly distributed (with respect to the Haar measure) random orthogonal matrix R was sampled. A vector of eigenvalues e was generated by sampling independent zero-mean unit-variance normal variates $n_i \sim \mathcal{N}(\cdot; 0, 1) \forall i \in \{1, ..., D_B\}$ and then setting $e_i = s_1 \tanh(s_2 n_i) \forall i \in \{1, ..., D_B\}$, with $s_1 = 6$ and $s_2 = 2$ in the experiments. This generates eigenvalues concentrated near $\pm s_1$ with this empirically observed to lead to systems which tended to be highly multimodal. A symmetric matrix $V = R \operatorname{diag}(e)R^T$ was then computed and the weights W set such that $W_{i,j} = V_{i,j} \forall i \neq j$ and $W_{i,i} = 0 \forall i \in \{1, ..., D_B\}$. The biases b where generated using $b_i \sim \mathcal{N}(\cdot; 0, 0.1^2) \forall i \in \{1, ..., D_B\}$. An example of a two-dimensional projection of independent samples from a Boltzmann machine relaxation density with D = 27 ($D_B = 28$), and W and b generated as just described in shown in figure 5. As can be seen even when projected down to two-dimensions the resulting density shows multiple separated modes.



Figure 5: Two-dimensional projection of 10000 independent samples from a Gaussian mixture relaxation of a Boltzmann machine distribution. The parameters W and b of the Boltzmann machine distribution where generated as described in Section C, with here $D_B = 28$ (rather than $D_B = 30$ as in the experiments) as independent sampling from larger systems exceeded the memory available on the workstation used. The two components shown correspond to the two eigenvectors of the generated basis R with the largest corresponding eigenvalues.

D Importance weighted autoencoder test images

84496008754346030877493753990998817782844828 974 22343080720401836449909945832 3 3 41 З 787 84687410076110726108539477116 0 5 63 2 99 7898 778179(472843578036849180803681488250 299 93 7010491009560034804077986746249040626161 31852047130292043978510183260049122 8936 1 1353102077647628365368801616414498443 398 614849671479218182197395409884441181948 323921986891603153401640813734456130238 ъ $\dot{\varphi}$ 262918908080098851651216634780881644915 1021836524319392870336194418595909989 3 86 173418161139707809370919809093386288341 6 32/360398234829498289367866332860219026 ô 8470642406174938/238499715980937797979629 96731767939909750992980886932939994805/87 1046140289637/223137919402342888260400912 90765681402239201010405302559948031)2071 022873040792010025040302259004446201618 242169089194019484281789880701374081348 9 1950 Ø033647806495889608757493103018991403E • 613361707685433956964611900513688124343 3239082965831200378437642144884182313182 3255069965435619368534042188481469446964

Figure 6: MNIST test images.

那台族路回是月后来 30时的医翻题条文作之母为文世国之间外国中局开始日本的运动的方式 空 开山石立》为功力或局势为古江郡,即端书外东下水伊门在大府国田下省南于市上进 出口自口中有所可益也所有自己問等自可自己自己在公司目前方了 (3)的目光读中和红田口 越先进死来又能力量的火艇固当了自己带口身子上来,整要也可被出火隐宅是不停了自己或重 上世武盛二年百姓的交易的成功的全部分别,我将华国日月日的日月日,在这时间上午上午前, 现代时间带 中门门间的 我们是这条出去自己走到下这样总人也站在个田町轮 又到新田西省上 Ş 物的是个可以拥有大品中国之际复展力的效器器器上包装出中台部的空空的外面的现在分词 1 医卵吸虫网膜发展网膜发生的可感感感的脑病性人的感情的 法正确的成长 医电子口疗术测虑 资料,只能加上了不能接近的外口的出口?但你可以知道我们不是不是不是不可能能 NEBWIE TO & BEAY WENDERNOOD NAT MO BELLE NELEDBURY 如此我会在避着了这些这里有了原来在那些在这些这些这些这些这些这些这些在这些在这些在这些。 的乌西尔姓代 建立电声力声的空机端推差着性敏气望温台氏后有脊索力 电变压压力地扩展负于 * 四日准督院确保(留了)和主任日子前的月日期所带着的现代和美国公共进行中国地方的日常 现住明 XD 》也能与国府同学部局最高量减多从大风限力公司地路军府国度过程及登出军市 就在到日月日前上的发出那个月中口,就不可以就要有些有些小户子从就必有了她,我们的了。" 5 鹅薄云发皮被皮掷上放用自住下的各面或火威破毁的那身后需有几个多肉还是当次为并让 日間有理は近代をというするののとうなないいのないでのないでものを見てきないのでの NEW BUERNAR ANNERSESES TO ARE SEE BAAR A GOODE . S FAX 机合压盖的的 的过去分词 法自己保留 医脊髓管膜 化化合物 化化合物 化化合物 化化合物

Figure 7: Omniglot test images.