

Annealed Stein variational gradient descent with localized Gaussian kernel

Zuheng(David) Xu

Student number: 24617185

ZUHENG.XU@STAT.UBC.CA

1. Summary of Stein variational gradient descent

Stein variational gradient descent (SVGD) (Liu and Wang, 2016) is a particle-based variational inference method, which iteratively transforms a set of samples from a reference distribution to the target distribution. In general, SVGD only requires the gradient of the log density function of the target distribution. This fits well into the framework of Bayesian inference, which usually has intractable normalizing constant in target distribution. Comparing to traditional sampling methods such as Markov Chain Monte Carlo (MCMC), SVGD uses deterministic updates and does not experience slow convergence due to autocorrelation between samples. A big advantage is that it leverage the gradient information of the target distribution, However, it also has many shortcomings. In this section, we briefly summarize the development of SVGD and its underlying theory, and discuss its limitations at the end of this section.

1.1 KL minimization and Stein's operator

Similar to normalizing flow (Kobyzev et al., 2020), SVGD aims to learn an invertible transformation from a reference distribution to the target distribution p . Instead of composing a series of parametric functions and optimizing over all the parameters, SVGD iteratively performs a non-parametric transformation to the current distribution and decreases the KL divergence to the target.

Let X be a random variable follows a distribution with density function $q(x)$, which dominates p . Let $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be an invertible and continuously differentiable map. In particular, we focus on the perturbed identity map:

$$T(x) := x + \epsilon \phi(x), \quad \phi \in \mathcal{H}_k^d \text{ and } \phi \text{ is smooth.}$$

Here \mathcal{H}_k^d is a *Reproducing kernel Hilbert space* (RKHS), a Hilbert space induced by a positive definite kernel $\kappa(x, x') : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$. We denote q_T as the push forward measure of $Y := T(X)$, of which the density function can be written as

$$q_T(y) = q(T^{-1}(y)) \cdot |\det \nabla T^{-1}(y)|.$$

The development of SVGD arise from a critical connection between the KL divergence and the *Stein's operator*.

Considering the directional derivative of $D_{\text{KL}}(q_T || p)$ along $\phi \in \mathcal{H}_k^d$,

$$\begin{aligned} \left. \frac{d}{d\epsilon} D_{\text{KL}}(q_T || p) \right|_{\epsilon=0} &= - \int q(x) (\nabla \log p(x)^T \phi(x) - \text{tr} \nabla \phi(x)) \\ &= -\mathbb{E}_q [\nabla \log p(X)^T \phi(X) + \text{tr} \nabla \phi(X)]. \end{aligned}$$

A key observation here is that $\nabla \log p(X)^T \phi(X) + \text{tr} \nabla \phi(X)$ is the *Stein operator* applied to $\phi(x)$. Then since $\phi \in \mathcal{H}_k^d$, the reproducing property of the kernel $\kappa(\cdot, \cdot)$ yields that

$$\left. \frac{d}{d\epsilon} D_{\text{KL}}(q || p) \right|_{\epsilon=0} = - \langle \phi, \mathbb{E}_q [\kappa(x, \cdot) \nabla \log p(X)^T \phi(X) + \nabla_x \kappa(x, \cdot)] \rangle_{\mathcal{H}_k^d}.$$

Note that the gradient is the directional derivative along the steepest descent direction. By maximizing the right hand side over $\phi \in \mathcal{H}_k^d$, we obtain that

$$\phi^*(\cdot) = \mathbb{E}_q [\kappa(x, \cdot) \nabla \log p(X)^T \phi(X) + \nabla \kappa(x, \cdot)], \quad \text{and} \quad (1)$$

$$\nabla_\epsilon D_{\text{KL}}(q||p)|_{\epsilon=0} = \|\phi^*\|_{\mathcal{H}_k^d}^2 = d_\kappa(q, p). \quad (2)$$

Here $d_\kappa(q, p)$ is the *kernelized Stein discrepancy* between q and p . Thus, by iteratively applying the following transformation to X_0 ,

$$X_0 \sim q_0, \quad X_{k+1} \leftarrow T(X_k), \quad T(X_k) := X_k + \epsilon \phi^*(X_k), \quad k = 0, 1, \dots \quad (3)$$

We are able to construct a non-parametric normalizing flow that decreases the KL-divergence sequentially. As revealed in Eq. (2), the fixed point of this transformation is achieved if and only if $d_\kappa(q, p) = 0$. It is worth pointing out that in general $d_\kappa(q, p) = 0$ does not imply that $q = p$. In fact, this requires a proper choice of the kernel function $\kappa(\cdot, \cdot)$, which will be discussed in detail in Section 1.2.

To implement this non-parametric normalizing flow Eq. (3), Liu and Wang (2016) approximate the optimal direction Eq. (1) using Monte Carlo samples. Specifically, one may draw n samples $\{x_i^{(0)}\}_{i=1}^n$ from the reference distribution p_0 and apply the empirical version of Eq. (3) to the set of particles, i.e., for $i = 0, 1, \dots, n$ and learning rate $\gamma_k > 0$, the numerical iteration at k -th iteration is as follows,

$$x_i^{(k+1)} \leftarrow x_i^{(k)} + \gamma_k \hat{\phi}_k^*(x_i^{(k)}), \quad \text{where} \quad (4)$$

$$\hat{\phi}_k^*(x) := \frac{1}{n} \sum_{j=1}^n \left[\kappa(x_j^{(k)}, x) \nabla \log p(x_j^{(k)}) + \nabla_{x_j^{(k)}} \kappa(x_j^{(k)}, x) \right].$$

Ideally, this iterative procedure will move the set of particle to the target distribution, which will eventually produce samples from p . The detailed description of the standard SVGD is presented in Algorithm 1. In practice, to achieve a faster numerical convergence, Liu and Wang (2016) suggest using the AdaGrad updates with momentum—also known as RMSprop—instead of vanilla gradient descent update described in Eq. (4). Specifically, RMSprop performs the update at k -th iteration as follows: for $i = 1, 2, \dots, n$,

$$g_{k+1} \leftarrow \eta g_k + (1 - \eta) \left(\hat{\phi}_k^*(x_i^{(k)}) \right)^2, \quad (5)$$

$$x_i^{(k+1)} \leftarrow x_i^{(k)} + \frac{\gamma_k}{\sqrt{g_{k+1} + \epsilon}} \hat{\phi}_k^*(x_i^{(k)}),$$

where $g_0 = 0 \cdot \mathbf{1}$ and all operations on vectors are elementwise. By default, the momentum factor η is set to 0.9 and the smoothing factor ϵ is set to 10^{-8} . Empirically, we do observe significantly faster convergence of the above update rule than the naive gradient descent update Eq. (4).

1.2 Choice of the kernel

Although we have formulated the construction of SVGD algorithm using some nice properties of RKHS, we have not provided any guideline of the selection of $\kappa(\cdot, \cdot)$. Note that the behaviour of SVGD is strongly dependent on the choice of the kernel. Therefore, it is helpful to understand the fundamental requirements of the kernel function.

Recall the transformation T described in Eq. (3), which is considered to be invertible and differentiable. By the Hadmard's global inverse function theorem, a sufficient condition for the invertibility of T is to have a bounded kernel function. Another concern on the kernel is related to the converging point of Eq. (3) as $k \rightarrow \infty$. As we mentioned above, the process of Eq. (3) will converge to a fixed point of the transformation T , where kernelized Stein discrepancy between the current

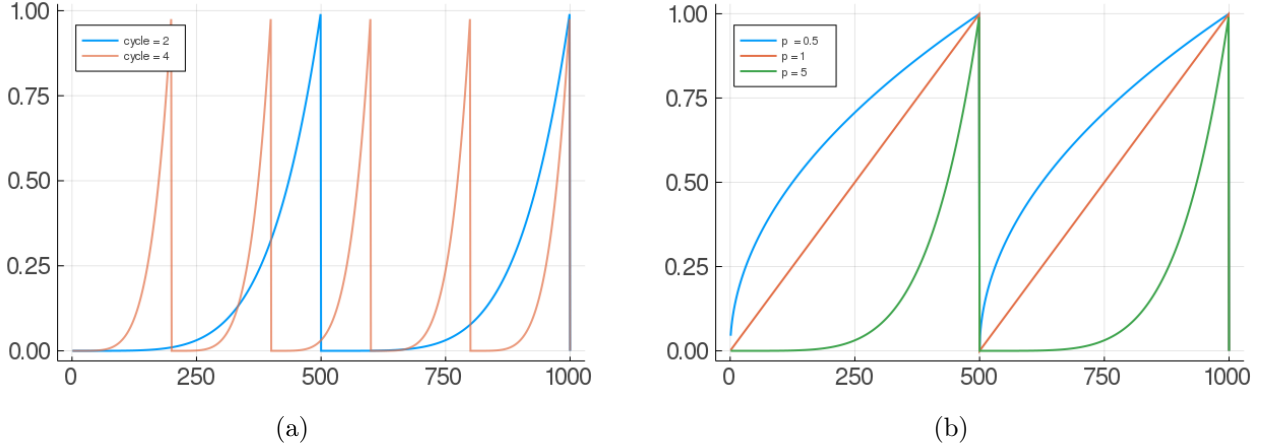


Figure 1: Annealing schedules. Fig. 1a shows two annealing schedules with 2 cycles and 4 cycles respectively. Fig. 1b shows annealing schedules across $p = 0.5, 1, 5$.

distribution q and the target distribution p arrives at 0. Ideally, we want to have $d_\kappa(q, p) = 0$ if and only if $q = p$. A sufficient condition for this is that $\kappa(\cdot, \cdot)$ is *integrally strictly positive definition* (Liu et al., 2016), i.e.,

$$\forall 0 < \|g\|_2^2 < \infty, \quad \int_{\mathcal{X}} g(x) \kappa(x, x') g(x') dx dx' > 0.$$

Examples of kernels satisfying these two properties are given by

$$\forall (x, x') \in \mathbb{R} \times \mathbb{R}, \quad \kappa(x, x'; h) = \exp\left(-\frac{1}{h}|x - x'|^p\right), \quad p \in (0, 2].$$

When $p = 2$, this yields the RBFkernel. Here $h > 0$ denotes the bandwidth of the kernel, which is a critical hyperparameter.

Different values of the kernel bandwidth can significantly influence the quality of SVGD; users need to hand tune the bandwidth to optimize the performance of SVGD. When using RBF kernel, Liu and Wang (2016) suggests picking the bandwidth based on the median of the pairwise Euclidean distances between the current particles, i.e., $h = \text{med}^2 / \log n$, where med is the median of the Euclidean distances. This procedure will update the bandwidth automatically at each iteration and works generally well. However, we will show in Section 2 when the median trick can be problematic.

1.3 Limitations of SVGD

A key limitation of SVGD is its computational burden. Note that the deterministic updates described in Eq. (4) require the evaluation of the kernel matrix and its gradient, which costs $\mathcal{O}(n^2)$. This makes SVGD less practical compared to other particle sampling methods such as sequential Monte Carlo sampler—which costs $\mathcal{O}(n)$ —if we aim to obtain a large number of samples from the target distribution.

The computation problem also occurs when calculating the gradient $\nabla \log p(x)$ for all particles; this is particularly expensive in a large dataset setting. Considering the target distribution as the posterior distribution of a Bayesian model, i.e., $p(\theta) \propto \pi_0(\theta) \prod_{i=1}^m p(x_i | \theta)$ where π_0 denotes the prior distribution of θ , computing $\nabla \log p(\theta)$ exactly is not practical as m is large. Stochastic gradient estimates using subsampled data points can be applied in this case to make SVGD scalable.

In the following sections, we focus on another common issue of SVGD when the target distribution $p(x)$ has multiple modes: the mixing problem, where the particles collapse at a particular mode of

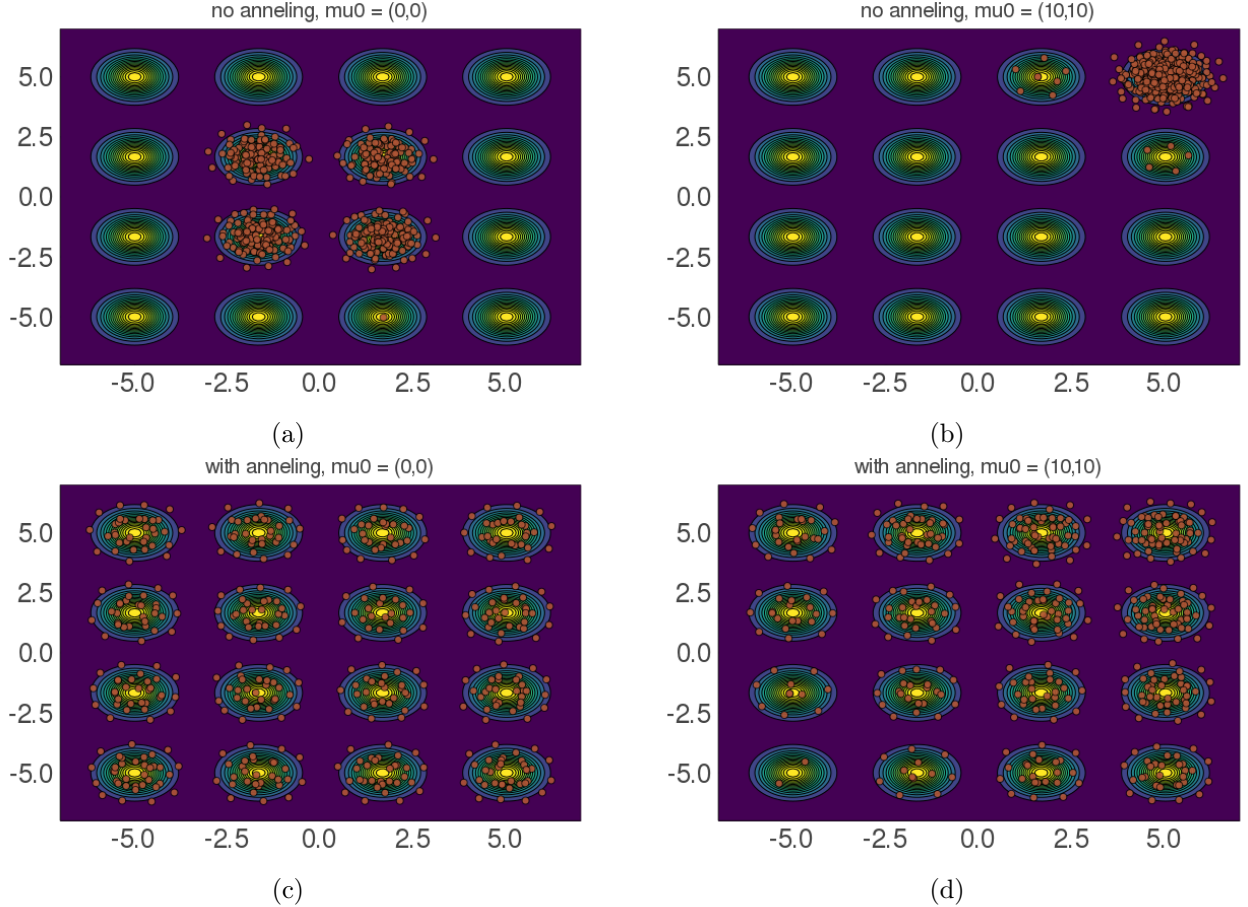


Figure 2: Comparison of SVGD (Figs. 2a and 2b) and A-SVGd (Figs. 2c and 2d) on the synthetic Gaussian mixture with 16 components with $\sigma = 0.5$ across two set of initial particles.

the target distribution and fail to distribute the particles to other high density regions. This is also known as the *mode collapse* phenomenon (D’Angelo and Fortuin, 2021; Zhuo et al., 2018) and appears often when the modes are distant from each other. The problem of mode collapsing tends to be more severe in high-dimensional settings, but can also be observed on a simple 2D Gaussian mixture target. Variants of SVGD are developed to resolve this issue (or at least alleviate the problem). For example, D’Angelo and Fortuin (2021) adopts the tempering method from the MCMC literature (Zhang et al., 2019) that includes an annealing schedule in the SVGD updates; Zhuo et al. (2018) leverages the latent structure of the statistical model and converts the inference problem into a lower-dimensional space.

In the next section, we provide a clear illustration of the mixing issue on two synthetic examples and introduce the annealed Stein variational gradient descent (A-SVGd), which tends to reduce the mixing issue. Further, we identify a central problem regarding the kernel bandwidth, which inspires our idea in Section 3.

2. Annealed Stein variational gradient descent

In this section, we provide a brief introduction of the annealed SVGD (A-SVGd), which is developed to address the mixing issue (D’Angelo and Fortuin, 2021), and use a few synthetic examples to identify a key limitation of the current work. The design of the synthetic example aims to understand the influence of the kernel bandwidth on the behaviour of SVGD or A-SVGd and help us to build the

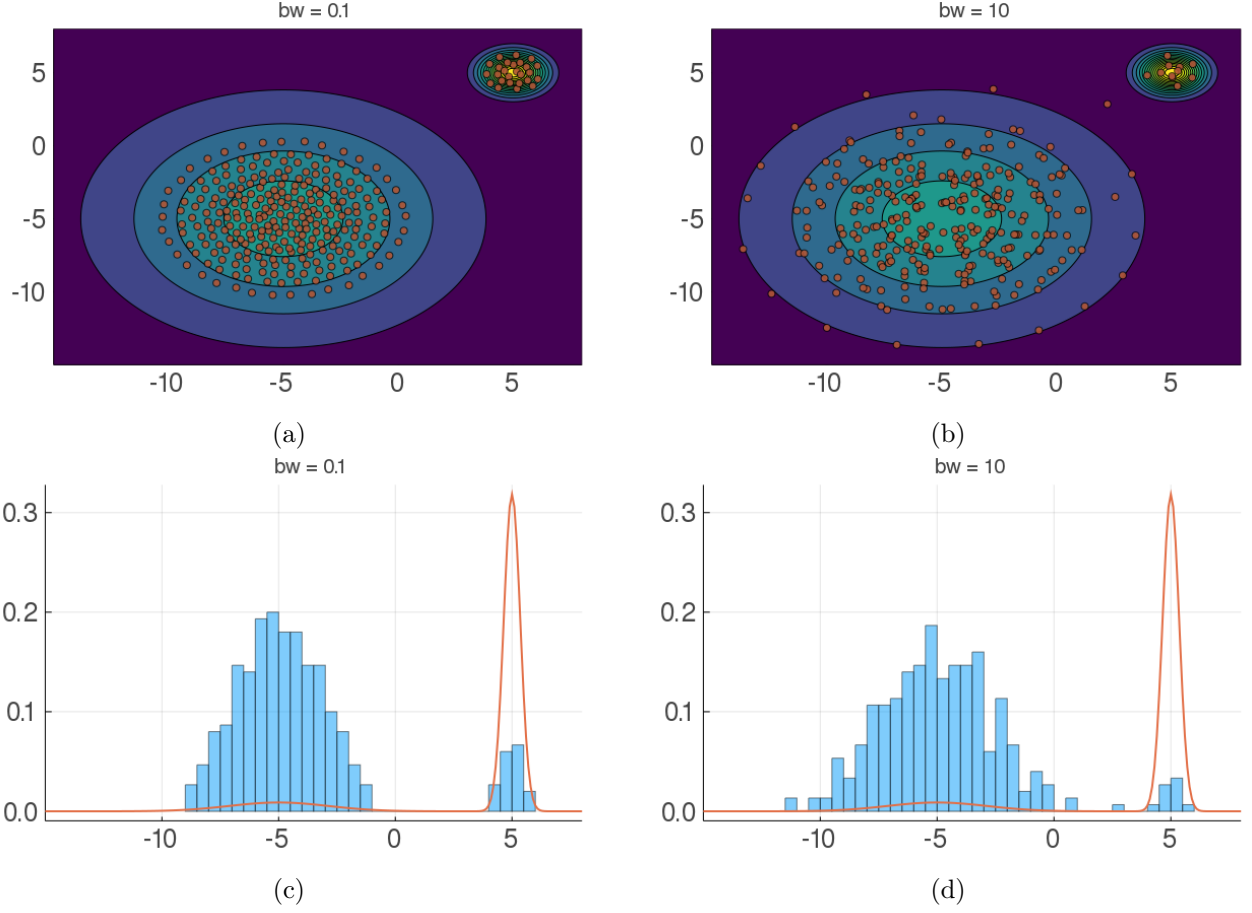


Figure 3: Comparison of A-SVGD on a 2-component Gaussian mixture with different bandwidth. Figs. 3a and 3b show scatter plots of particles at the last iteration. Figs. 3c and 3d show histograms of the particles by projecting the locations to the line $y = x$, where the orange curve denotes the density of Gaussian mixture along the slice $y = x$.

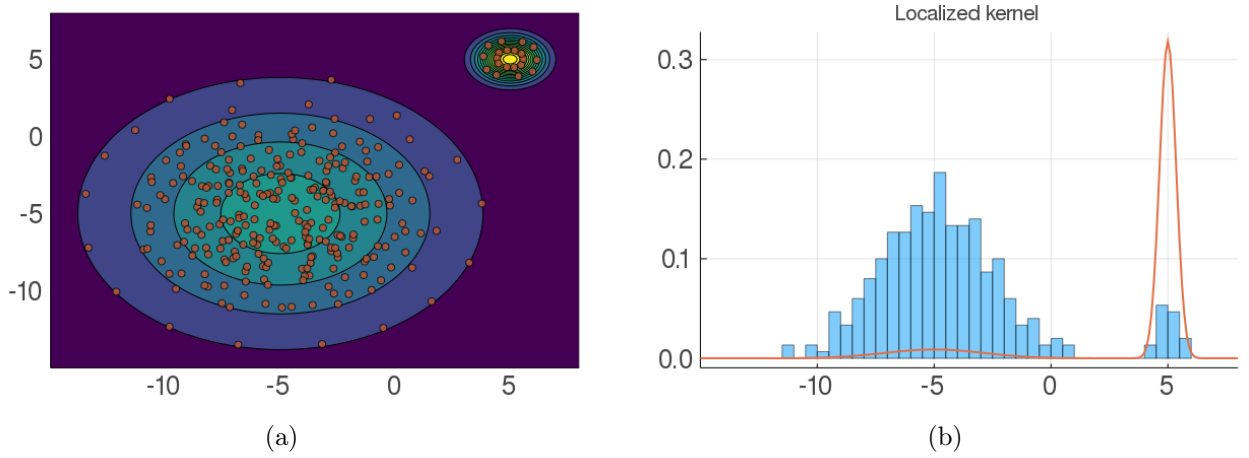


Figure 4: A-SVGD with adaptive local bandwidth on a 2-component Gaussian mixture.

intuition of picking a proper value. All the experiments performed in this report use the RMSprop

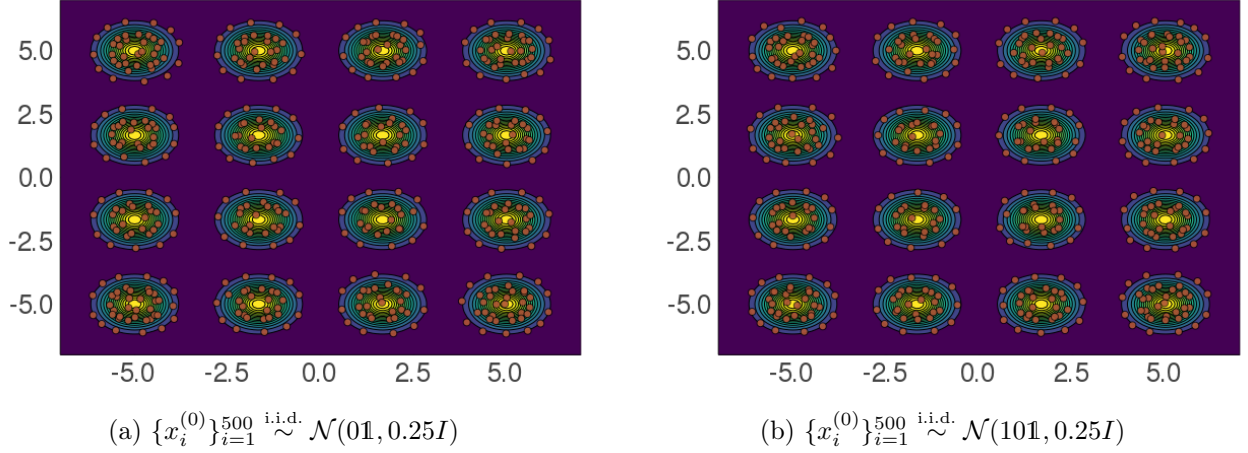


Figure 5: A-SVGd with adaptive local bandwidth on a 4×4 Gaussian mixture.

update rule Eq. (5). We run both SVGd and A-SVGd with standard RBF kernel for 1000 iterations under a constant learning rate $\gamma_k = 0.1$.

We start with the development of A-SVGd. Considering the updates Eq. (4), the optimal descent direction $\hat{\phi}^*(x)$ for the particle x includes two parts: the driving force that moves the particle to a high density region of $p(x)$, and a repulsion that prevents the particles collapsing together. The specific derivation is as follows,

$$\hat{\phi}_k^*(x) = \frac{1}{n} \sum_{j=1}^n \left[\underbrace{\kappa(x_j^{(k)}, x) \nabla \log p(x_j^{(k)})}_{\text{driving force}} + \underbrace{\nabla_{x_j^{(k)}} \kappa(x_j^{(k)}, x)}_{\text{repulsion}} \right].$$

When mode collapse occurs, particles produced by SVGd are trapped at some region and fail to explore the full support of $p(x)$. This often happens when the driving force dominates the repulsion. An intuitive solution is to boost the repulsion or scale down the driving force to enable a wider spread of the particles. D’Angelo and Fortuin (2021) follows this intuition and modify the descent direction by introducing an annealing term $\alpha(k) \in [0, 1]$:

$$\hat{\phi}_{\alpha_k}(x) = \frac{1}{n} \sum_{j=1}^n \left[\alpha(k) \kappa(x_j^{(k)}, x) \nabla \log p(x_j^{(k)}) + \nabla_{x_j^{(k)}} \kappa(x_j^{(k)}, x) \right]. \quad (6)$$

In particular, the authors suggest using a cyclical annealing schedule, i.e.,

$$\alpha(k) = \left(\frac{\text{mod}(k, K/C)}{K/C} \right)^p,$$

where K is the total number of iterations, C is the number of cycles, and $p > 0$ is an exponent factor that controls the speed of the transition between two phases. As demonstrated in Fig. 1, smaller p corresponds to a shorter period of the exploratory phase. In all our experiments using A-SVGd, we set $p = 1$ and $C = 2$. **add a few sentence about cyclical**

Now we compare the performance of A-SVGd and SVGd on a synthetic Gaussian mixture (Fig. 2). The target distribution is a 2D mixture of 16 Gaussian component with means equally distributed on a 4×4 grid with identical covariance $0.25I$. We consider two different initialization particles: Figs. 2a and 2c are initialized with 500 i.i.d. samples from $\mathcal{N}(0, 0.25I)$; Figs. 2b and 2d are initialized with 500 i.i.d. samples from $\mathcal{N}(10, 0.25I)$. In this example, we use RBF kernel with $h = 0.5$. Note that in both initialization schemes, A-SVGd is able to explore all the modes while all the particles produced by

SVGD collapse at the components close by the initialization. Comparing Figs. 2c and 2d, one may notice that the performance of A-SVGD still depend on the location of the initial particles.

In this experiment, we create a pathological yet simple example such that neither small nor large bandwidth is satisfying. Considering a two-component Gaussian mixture:

$$\frac{1}{2}\mathcal{N}(-5\mathbf{1}, 9I) + \frac{1}{2}\mathcal{N}(5\mathbf{1}, 0.25I).$$

Fig. 3 display the results of A-SVGD with 300 initial particles sampled from $\mathcal{N}(0\mathbf{1}, 25I)$ using different kernel bandwidth $h = 0.1, 10$. With small h particles located in the wide Gaussian component are overly contracted; with large h , almost all particles converges to the wide mode. This is essentially due to the fact that the curvature of the local basin of two gaussian components is quite different. When using a wide kernel bandwidth, the driving force of those particles around the narrow Gaussian component is significantly influenced by gradient information $\log p(x)$ of the particles that are far away. As most particles are trapped in the wider Gaussian component, the driving force tends to point to the wide mode. Also, a large bandwidth leads to a stronger repulsion, which encourages the spread of particles. On the contrary, with small bandwidth—as shown in Figs. 3a and 3c—although the peaky Gaussian absorbs more particles, due to a weak repulsion introduced by the small bandwidth value, particles at the wider Gaussian component are not spread enough, resulting in an underestimation of the variance. This hints at an important fact that using a universal kernel bandwidth for all particles is not ideal, and we should use different kernel bandwidths when updating each of the particles depending on its local basin of $\nabla \log p(x)$. In the next section, we propose an adaptive local bandwidth selection scheme that requires no input from the user.

3. A-SVGD with adaptive localized kernel

With the intuition obtained in the last section, we suggest using different bandwidth values for each particle and updates the bandwidth adaptively across the iteration. The goal of this section is to design an adaptive procedure that learns a proper kernel bandwidth for each particle automatically. As illustrated in Fig. 3, a careful choice of kernel bandwidth is necessary when the target distribution has multiple modes and one might consider using different bandwidth values at a different location. Intuitively, a proper choice of the bandwidth should reflect the local geometry of $\log p(x)$ —the curvature of $\log p(x)$ —which includes information of the Hessian $\nabla^2 \log p(x)$.

We consider the anisotropical Gaussian kernel with covariance Σ ,

$$\kappa(x, x'; \Sigma) = \exp(-(x - x')^T \Sigma^{-1} (x - x')).$$

For computational reason, we propose to set Σ to the absolute value of $\text{diag}(\nabla^2 \log p(x))$.

Based on the empirical findings in Section 2, we find narrower mode requires smaller bandwidth while wider mode prefers a larger kernel bandwidth. Therefore, when updating x_i , an ideal choice of h would be $1/\sigma_1(\nabla^2 \log p(x_i))$, where $\sigma_1(M)$ denotes the maximal singular value of matrix M . However, solving $\sigma_1(\nabla^2 \log p(x_i))$ is generally too expensive. We then use the maximal diagonal element of $\nabla^2 \log p(x_i)$ as approximation. The complete procedure of A-SVGD using adaptive local bandwidth is described in Algorithm 2.

With other settings fixed, we apply Algorithm 2 to the previous two synthetic examples. As demonstrated in Fig. 4, by using the local bandwidth, the samples in the wide component correctly characterize its variance and there is a certain amount of samples concentrate at the peak mode. This combines the advantages of both small h and large h . And for the 16-component Gaussian example, since all Gaussian component shares the same covariance, we do not expect A-SVGD with adaptive local bandwidth could outperform the previous results significantly. As shown in Fig. 5, the result is quite similar to the A-SVGD using hand-tuned bandwidth; but importantly this requires no user input in terms of kernel bandwidth.

It is worth pointing out that leveraging the Hessian of $\log p(x)$ introduces additional computation complexity. Evaluating the Hessian for all particles costs $\mathcal{O}(d^2N)$, denoting a quadratic growth as d increases. This makes SVGD less practical in a high dimensional setting. In that case, we might consider using Gauss-Newton approximation to the Hessian matrix or other methods that approximate the second-order derivative using first-order information. Also, when n is large, evaluating $\nabla^2 \log p(x)$ for all the particles is impossible. But employing stochastic estimates of the exact Hessian matrix is a reasonable approach.

4. Experiment

5. Discussion

some open problems: post hoc evaluation of the mixing issue.

is there optimal annealing schedule.

is there any geometric interpretation of bandwidth

References

- D’Angelo, F. and Fortuin, V. (2021). Annealed Stein variational gradient descent. *arXiv preprint arXiv:2101.09815*.
- Kobyzev, I., Prince, S., and Brubaker, M. (2020). Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Liu, Q., Lee, J., and Jordan, M. (2016). A kernelized stein discrepancy for goodness-of-fit tests. In *International Conference on Machine Learning*, pages 276–284.
- Liu, Q. and Wang, D. (2016). Stein variational gradient descent as gradient flow. In *Advances in Neural Information Processing Systems*, pages 2370–2378.
- Zhang, R., Li, C., Zhang, J., Chen, C., and Wilson, A. (2019). Cyclical stochastic gradient mcmc for bayesian deep learning. *arXiv preprint arXiv:1902.03932*.
- Zhuo, J., Liu, C., Shi, J., Zhu, J., Chen, N., and Zhang, B. (2018). Message passing Stein variational gradient descent. In *International Conference on Machine Learning*, pages 6018–6027.

Appendix A. Local bandwidth of two synthetic example

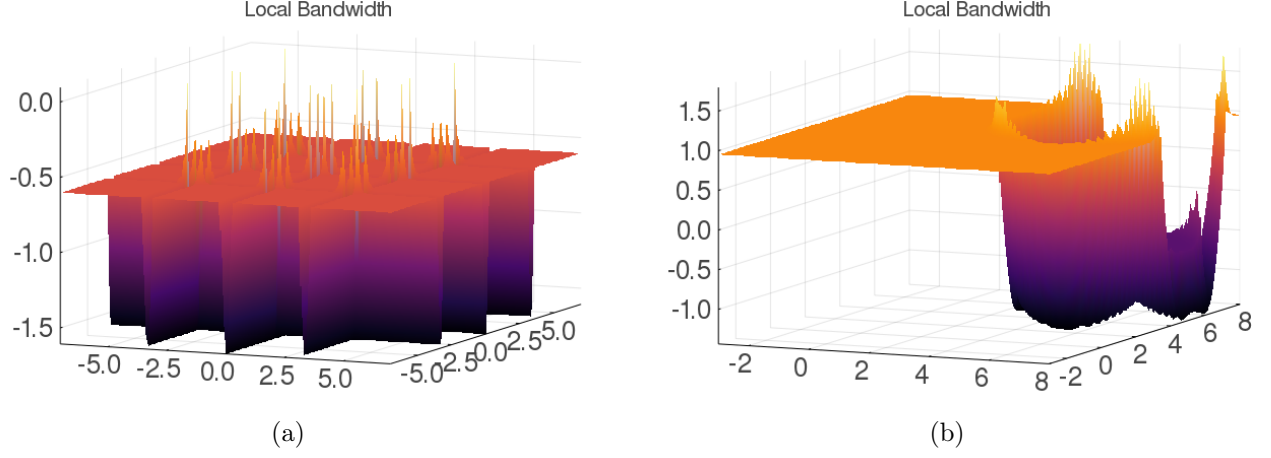


Figure 6: Local bandwidth value for 4×4 Gaussian grid (Fig. 6a) and 2-component Gaussian mixture (Fig. 6b).

Appendix B. Algorithms

Algorithm 1 SVGD

```

procedure SVGD( $\{x_i^{(0)}\}_{i=1}^n, \nabla \log p(x), \kappa(\cdot, \cdot), \gamma_k$ )
  for  $k = 0, 1, \dots, K - 1$  do
    for  $i = 1, 2, \dots, n$  do
       $x_i^{(k+1)} \leftarrow x_i^{(k)} + \gamma_k \hat{\phi}_k^*(x_i^{(k)})$  where  $\hat{\phi}_k^*(\cdot)$  is defined in Eq. (4)
    end for
  end for
  return  $\{x_i^{(K)}\}_{i=1}^n$ 
end procedure

```

Algorithm 2 A-SVGd with local bandwidth

```

procedure SVGD( $\{x_i^{(0)}\}_{i=1}^n, \nabla \log p(x), \nabla^2 \log p(x), \gamma_k, \alpha(t)$ )
  for  $k = 0, 1, \dots, K - 1$  do
    for  $i = 1, 2, \dots, n$  do
       $\Sigma_i \leftarrow 1 / (\max \{ \text{diag } \nabla^2 \log p(x_i^{(k)}) \})$ 
       $\kappa(x, x'; \Sigma_i) \leftarrow \exp(-(x - x')^T \Sigma_i^{-1} (x - x'))$ 
       $x_i^{(k+1)} \leftarrow x_i^{(k)} + \gamma_k \hat{\phi}_{\alpha_k}(x_i^{(k)})$  where  $\hat{\phi}_{\alpha_k}(\cdot)$  is defined in Eq. (6)
    end for
  end for
  return  $\{x_i^{(K)}\}_{i=1}^n$ 
end procedure

```
