

1 **[Reviewer #1]** Thank you for your valuable comments and suggestions.

2 *On choosing a matrix kernel in practice:* As the reviewer correctly noted, the problem of choosing an optimal kernel is
3 a very hard problem. The difficulty lies on the fact that ideal optimal kernel should incorporate the intrinsic geometric
4 information, but leads to no expensive computation.

5 In this work, we propose the mixture preconditioning kernel as a practical heuristic to strike the balance between
6 efficiency and computational cost. Compared with the constant preconditioning (which corresponds to using constant
7 Hessian matrix), the mixture preconditioning allows us to incorporate different curvatures in different locations, which
8 is of critical importance for complex models such as the "star" in Figure 1. Meanwhile, compared with the more
9 theoretically motivated kernel derived with change of variables, the mixture preconditioning kernel avoids taking high
10 order derivatives and is much more computationally efficiently.

11 **[Reviewer #2]** Thank you for your valuable comments. We will improve our presentation based on your questions.

12 *Comparisons with SVNM :* The advantage of matrix SVGD over SVN is clear both theoretically and empirically. As we
13 point out in L188-199, SVN can be in fact viewed as solving the same fixed point equation of SVGD with a scalar
14 kernel and hence can not leverage the advantage of matrix kernels. For empirical comparison on the realistic examples,
15 we did test the official Matlab code provide by the SVN authors, but the results we obtained were much worse than our
16 methods and baselines (similar to the results in Fig 1), and hence did not investigate it further. For example, the RMSE
17 score of SVN of BNN on the Boston dataset is 3.1830(± 0.1829) (averaged on 20 random trials), which is much worse
18 than all the methods reported in Table 1 (< 2.9). We will investigate further and add discussion on this issue.

19 *The unit ball is not unique:* The RKHS $\mathcal{H}_k^d := \mathcal{H}_k \times \dots \times \mathcal{H}_k$ must be equipped with inner product $\langle f, g \rangle_{\mathcal{H}_k^d} =$
20 $\sum_{i=1}^d \langle f_i, g_i \rangle_{\mathcal{H}_k}$, which corresponds to using L_2 norm on \mathbb{R}^d . Using other norms on \mathbb{R}^d can not define an RKHS.

21 *Is R-SVGD a special case of our method:* No, we will acknowledge out this.

22 *[CBBGGMO2019] :* Thanks for pointing it out. We were not aware this work by the time of submission. We will add a
23 reference. But note that the preconditioned kernel in this paper is still a scalar kernel and hence substantially different
24 from our matrix kernel.

25 *Preconditioning matrices:* We use Hessian matrix in Sec 4.1 and Fisher with K-FAC approximation in Sec 4.2-4.4.

26 *How many particles:* For the toy example, we use 50 particles for all the cases. For other applications, please see *line
27 235*, *line 245* and *line 264* for more details.

28 *How many anchor points :* Please see *line 212* (we use all the particles as the anchor points).

29 *What kernel is used for computing MMD in Fig 1:* We use the standard RBF kernel for evaluating MMD.

30 *Different initialization?:* We did use the same initialization for all methods with the same random seeds. Please see Fig
31 5, 6, 7; or code/bnn/trainer.py:L293, or code/rnn/acl_sentence_classification.py: L384; in Figure 1(f), the plot is started
32 from the 10-th iteration to avoid large values for better visualization.

33 *Section 4.3 and 4.4: mini-batch size and #iteration :* For Section 4.3, we use a mini-batch size of 100 for training; for
34 large dataset such as year, we set the batch size to be 1000. In general, we grid search the least training epochs required
35 for our implementation of SVGD to have a performance that matches the results reported in the original SVGD paper;
36 we then train all other methods with the same training epochs. For Section 4.4, we use a mini batch-size of 50. we train
37 all algorithms for 20 epochs;

38 *the interchange of expectation and inner product should be justified:* It is true once integrations involved are finite.

39 *Reference:* We will add a reference on the median trick and Lemma 2.

40 **[Reviewer #3]** Thank you for your valuable comments and suggestions.

41 *L102: There looks to be a typo:* Yes, \mathcal{H}_k^d should be corrected to be \mathcal{H}_K here.

42 *Experiments 1 & 2, is there any way to compute the algorithms on CPU time rather than iterations :* As the other cases
43 when Newton or natural gradient are used, the timing comparison heavily depends on implementation, approximation
44 method of the matrix inversion, and many other factors such as matlab or python code and gpu-based or cpu-based
45 implementation. Another important point, unique to our method, is that the SVGD with matrix kernels can yield a set of
46 better fixed points than SVGD with scalar kernels, which may worth the effort even when the time-wise convergence is
47 slower (in comparison, the standard Newton's method and natural gradient yield the same set of fixed point as standard
48 gradient descent).