

1 Thanks to all reviewers for the insightful comments! We have already open-sourced most of our code (link hidden to
2 preserve anonymity) and will provide the link at the beginning of Section 5 of the final paper. Detailed responses below:
3 **R1: The overall method seems to be not end-to-end.** Algorithm 1 is in fact an end-to-end algorithm. All the trainable
4 parameters β and $\{W_l, b_l\}_{l \leq L-1}$ are updated simultaneously in the main SGD step. The only parameters that are not
5 updated are the W_L and b_L , which are fixed random features that are frozen throughout the learning. The precision
6 matrix update is only a side computation for summary statistics derived from the SGD-updated model. As a result, the
7 optimization procedure is very similar to the deterministic SGD training. We use the same initialization for the hidden
8 weights as in the deterministic model, and use the default Glorot uniform initializer for the GP output layer weights β .
9 **R1: Clarity regarding Eqn (5).** Thanks for highlighting a point that merits more clarification. Briefly, different from
10 a classic minimax problem, we derived (5) under the assumption that we have partial knowledge about p^* (i.e., we
11 know the domain probability $p^*(x \in X)$), therefore it is possible for the known property of p^* to enter into the final
12 expression. Please see Appendix B for a full statement of the motivation, the assumptions and the proof for Eqn. (5). In
13 particular, please see line 619-626 for an explanation of why (5) is structured as such. Eqn (5) (corresponds to Appendix
14 Proposition 2) is used to motivate Section 2.2. In the final paper, we will include additional explanations / pointers to
15 Appendix B around line 100, and replace "Proposition 2" with "Equation (5)" on line 110 to improve clarity.
16 **R2: It seems that the \mathcal{X} is with L_2 -norm as distance metric.** We in fact allow the metric for \mathcal{X} to be non-Euclidean
17 so it reflects the semantically meaningful distance in the data space (please see statement on line 55-56, the discussion
18 on line 141-149, and the proof for Proposition 1 in Appendix D.2. which does not impose restriction on $\|\cdot\|_{\mathcal{X}}$ and
19 is based on the theoretical work of [3]). In addition, we'd like to point out that in the vision / language experiments
20 (Section 5.2), the SNGP has superior performance in distinguishing in-domain / out-of-domain data, which is not likely
21 if SNGP can only preserve a L_2 metric, which is not suitable for an image / language manifold.
22 **R3: Ablation study.** We have conducted such ablation study in Appendix C, where DNN-SN and DNN-GP are ablated
23 versions of SNGP. Figure 2-3 shows that in the 2D example, the uncertainty surface of a DNN-SN behaves similarly to
24 a deterministic DNN, while that of a DNN-GP is lacking in preserving input distance. Table 4-6 shows in the vision and
25 language experiments, DNN-SN and DNN-GP tend to outperform the deterministic baseline, but underperform SNGP.
26 **R3: Comparison to methods designed explicitly for OOD detection.** Please see table for performance compar-
27 ison to popular OOD methods evaluated using area under precision recall curve (AUPRC) (we will add it to
28 Appendix C). We denoted CIFAR-10/-100 as C10/100. As shown, despite not designed explicitly for OOD, SNGP is
29 competitive and sometimes outperforms other OOD approaches, especially on difficult near-OOD tasks (e.g, CIFAR 10
30 v.s. 100 and vice versa). Mahalanobis = Mahalanobis with feature ensemble and input processing.
31 **R4: Unclear why the distance to the training data should be used for the uncertainty measure / whether this
32 distance awareness property is indeed advantageous.** Intuitively, given a testing example that "looks different" from
33 the training data (i.e., far from the training data manifold), a model's uncertainty metric is expected to return a high value
34 (see, e.g., Fig 1a). Such definition of model uncertainty (or "epistemic" uncertainty) in terms of distance/dissimilarity
35 from observed data has been widely adopted in both the UQ and the ML literature (c.f. Kiureghian and Ditlevsen
36 (2009).Aleatory or epistemic? Does it matter? , Kendall and Gal (2017).What Uncertainties Do We Need in BDL for CV?, *NeurIPS* and the many papers
37 citing them), and empirically can be measured by a model's OOD accuracy (see Table 1-2, and Table 4-6 in Appendix).
38 Quoting other reviewers: "Empirical results on benchmark datasets show superior performance in OOD." (R1), "This
39 work conducted convincing experiments on various datasets...showed the advantage of the proposed method." (R2), etc.
40 **R4: Why the special setting in this paper is inevitable...Either the theoretical derivations or the empirical results
41 do not support why the practitioners have to use the proposed modification.** Contrary to local methods, vanilla
42 DNN models tend to have difficulty in achieving the distance-preservation property shown in Eqn (6). For example,
43 vanilla DNNs are found to be vulnerable to adversarial examples - they can be sensitive to tiny perturbation in the input
44 space, yet sometimes insensitive to semantics-altering edits to the training data - i.e., not input distance aware [33,34].
45 To this end, the paper's theoretical result (Proposition 1) ensures SNGP's ability in guaranteeing distance preservation,
46 and the empirical result (Table 1-2, and Appendix C.2) shows that such modification leads to concrete improvement in
47 ECE/OOD performance when compared to an unmodified baseline.
48 **R4: ...disentanglement...is contrary to the explanation that the distance-preservation matters.** Disentanglement
49 and distance-preservation (i.e. invariance) are both important properties for a representational learning algorithm, and
50 they do not contradict each other (see, e.g., Achille and Soatto (2018).Emergence of Invariance and Disentanglement in Deep Representations, *JMLR*).
51 For a DNN hidden mapping $h : \mathcal{X} \rightarrow \mathbb{R}^d$ which is a coordinate transform from the input space to a hidden space
52 $h(x) \in \mathbb{R}^d$, *disentanglement* describes $h(x)$'s ability in separating salient latent features from the noise among its
53 d coordinates, while *distance preservation* describes $h(x)$'s ability in translating a semantically meaningful (often
54 non-Euclidean) measure in the data manifold into that in the Euclidean space [30]. With suitable model specification,
55 disentanglement can happen jointly with distance preservation (e.g., see Figure 2 of [30]). There have been many work
56 that try to achieve both for the purpose of generalization and adversarial robustness, notably via Lipschitz regularization
57 or invertible (i.e. bi-Lipschitz) networks [35, 69] (also, e.g., Engstrom. (2019) Adversarial Robustness as a Prior for Learned Representations)

Method / AUPRC	C10 vs SVHN	C10 vs C100	C100 vs SVHN	C100 vs C10
MSP+OE	89.4	76.2	52.9	32.6
Mahalanobis	99.1	-	98.4	-
ODIN	92.5	-	93.9	-
SNGP	99.0	90.5	92.3	80.1