

## 1                    **NeurIPS Rebuttal for “Reliable training and estimation of variance networks”**

2 We thank the reviewers for their constructive and fair reviews. We here address the key concerns, and note that the  
3 paper will be updated accordingly. We will discuss two shared concerns, and then move to individual reviewers.

4 First, however, we emphasize that the paper provides the first study of training issues related to variance networks, and  
5 develops new training methodologies that result in significant gains over current state of the art.

6 **High dimensionality (R2+R3):** Two reviewers expresses concern about extending our methodologies to higher  
7 dimensions. This is reasonable as both the locality sampler and the variance extrapolation depend on Euclidean  
8 distances, which are known to perform poorly in high dimensions. While we acknowledge the issue, we want bring  
9 additional perspectives to the discussion:

- 10        • One definition of variance is that it is the mean squared deviation from the mean, so the notion of variance  
11        itself is dependent on the Euclidean distance. This may indicate that variance networks are inherently difficult  
12        in high dimensions, where some regularity assumptions may be needed.
- 13        • While Euclidean distances behave poorly in high dimensions, there is plenty of empirical evidence that nearest  
14        neighbors are somewhat better behaved. As an example, spectral learning techniques have been very successful  
15        in moderately high dimensions. This is also in line with our observation that the proposed methods work well  
16        on MNIST, which is 768 dimensional (medium-sized data).
- 17        • As R3 mentions, the issue can possibly be circumvented through a different choice of metric; possibly one that  
18        is learned. This is an interesting venue for further research.

19 **Computational costs (R2+R3):** We agree that performing nearest neighbor search in the inner training loop presents a  
20 scalability issue. Two practical comments are, however, in order:

- 21        • We can rely on fast approximate nearest neighbor algorithms. This is an established area of research, where, in  
22        particular, randomized algorithms are currently showing notable speed gains.
- 23        • The variance network is typically trained on a GPU while the nearest neighbor search is performed on the  
24        CPU. Thus, the two processes can in principle be performed in parallel. The total training time can then be  
25        unaffected as long as the neighbor search is faster than a forward and a backward pass.

26 That being said, we acknowledge the concern, and view this as a topic for future work.

27 **We consider** the above two concerns to be both valid and important, but we also maintain that the paper provides the  
28 first working implementation of variance network training, which is a significant contribution.

29 **R1:** We agree with the reviewer that the three variance problems discussed in this paper (underestimation, trivialization,  
30 extrapolation) could easily each be the subject of its own paper. With this paper we aim at bringing attention to the  
31 multiple separate problems within uncertainty estimation, and argue that each problem requires its own solution. We  
32 provide initial solutions that significantly improve on the current state of affairs, but do expect that they can all be  
33 improved. We, further, would not be surprised if different application areas might benefit from slightly different  
34 solutions. We will update the paper with a discussion of these matters in more detail.

35 **R3:** We agree that the proposed methodologies do not solve the problem of overfitting. Interestingly, this problem also  
36 depends on the flexibility of the mean function. If this perfectly fits the data, then overfitting of the variance becomes  
37 more likely. This indicates that the link between variance estimation and regularization of the mean, could use further  
38 study; we consider this to be an interesting venue for research.

39 In practice we do not believe that the bias introduced by the sampling scheme plays much of a role, if the data is dense.  
40 However, theoretically we could have regions of data that are under-trained because they are underrepresented in the  
41 samples. The reason for the improved result is more likely due to an optimization that avoids the local minima of the  
42 usual mini-batching.

43 Regarding, Eq. 5: Yes, this is a typo. Thanks for catching it!