

## 1 **Response to Reviewer 1**

2 We sincerely thank Reviewer 1 for referring us to four relevant papers [1-4]. We will include a detailed discussion on  
3 these papers in our revised manuscript, but would like to emphasize that our work differs from, and is complementary  
4 to this literature. This is further explained below:

5 Paper [1] provides a very interesting relationship between Fisher divergence and Stein's operator, whereby Hyvarinen  
6 score and its extensions are cast as specific Stein's operators. Paper [4] establishes a more general result than that of S.  
7 Lyu.

8 In contrast, our work focuses on utilizing Hyvarinen score (and its extensions) to develop information quantities for both  
9 interpretability and computational benefits (e.g. in our example of fast tree approximation). That said the computational  
10 motivation still remains to be the getting rid of normalizing constants.

11 Regarding your first technical question: It follows from an application of weak law of large number that by minimizing  
12 the Hyvarinen score the parametric density converges to the shadow of the true data generating distribution under mild  
13 assumptions. If the model class is well-specified then convergence to the data generating distribution is guaranteed.

14 Regarding your second technical question raised: The Fisher entropy is intimately related to the Cramer-Rao Bound in  
15 parameter estimation. Fano's inequality also gives lower bounds of model selection/message decoding error (so larger  
16 bound implies more complexity for 'description'). We understand that Reviewer 1 may have a different interpretation  
17 of our wording, and we will clarify this in the revision.

18 [1] Liu, Qiang, Jason Lee, and Michael Jordan. 'A kernelized Stein discrepancy for goodness-of-fit tests.' ICML 2016.

19 [2] Chwialkowski, Kacper, Heiko Strathmann, and Arthur Gretton. 'A kernel test of goodness of fit.' ICML 2016.

20 [3] Gorham, Jackson, and Lester Mackey. 'Measuring sample quality with kernels.' ICML 2017.

21 [4] Liu, Qiang, and Dilin Wang. 'Stein variational gradient descent: A general purpose bayesian inference algorithm.'  
22 NIPS 2016.

## 23 **Response to Reviewer 2**

24 We appreciate Reviewer 2's comments and recommendations.

25 We will do another proof reading and remove typos. Due to page limit we could not include more extensive experiments,  
26 but we plan to do so in subsequent works.

27 Reviewer 2 raised an excellent point about the extension from pairs of variables to groups of variables. We have some  
28 recent results on this (not reported due to page limit), and hope to expand on them and present them in future work.

29 Regarding the two questions raised by the Reviewer 2, the values of  $p$  ranges from 2 to 11 and the joint density is  
30 determined by the conditional exponential family distribution in (12) and the discovered tree structure; the second term  
31 in (2) is  $q(y)$ , and  $p(y)$  has been implicitly included in the expectation.

## 32 **Response to Reviewer 3**

33 We appreciate Reviewer 3's comments and recommendations.

34 Reviewer 3 raised an excellent point about defining cross entropy. We think it is doable and it would be another analogy  
35 to the classical cross entropy. Despite theoretical interest, we have not yet found a real data application for this. We will  
36 include a discussion of this in the revision.

37 Like the point made in (Marsh 2013), it is also true in our case that the relative entropy between distributions is more  
38 interpretable than entropy on its own. We have not found an application of gradient mutual information in the settings  
39 of information bottleneck. But we believe this direction is extremely interesting. As suggested by Reviewer 3, we will  
40 do another proof reading and remove the typos.