

1 We thank the reviewers for their valuable and positive feedback. We will address their main concerns:

2 **Reviewer #1.** We thank the reviewer for the positive and encouraging feedback and address their suggestions and
3 questions: 1) *On Clarity “TPR and FDR same Figure”*: For space restriction, and in line with prior papers eg HRT [8],
4 we elected to have them in the same plot, but we will make sure to remind the reader of the expected trends in the legend
5 of the figure. 2) *Comparison to the Gate Formulation in (P)*: In our early experiments using a NN critic, we found
6 this formulation to be under-performing. Theoretically, Formulation (P) does not lead to a convex formulation in the
7 RKHS. (P) was indeed also studied in the cited reference [10] (Appendix A), using the Hilbert-Schmidt Independence
8 Criterion and a concave approximation of ℓ_0 , and was found to be under-performing as well. 3) *“Analysis of the
9 Consistency”*: This is an interesting question. In order to recover the correct conditional independence we elected to
10 use FDR control techniques to perform those dependent hypotheses testing (btw coordinates). By combining SIC with
11 HRT and knockoffs we can guarantee that the correct dependency is recovered while the FDR is under control. For the
12 consistency of SIC in the classical sense, one needs to analyze the solution of SIC, when the critic is not constrained to
13 belong to an RKHS. This can be done by studying the solution of the equivalent PDE corresponding to this problem
14 (which is challenging, but we think it can be also managed through the η -trick). Then one would proceed by finding
15 1) conditions under which this solution exists in the RKHS 2) Generalization bounds from samples to the population
16 solution in the RKHS. We can mention these points in the final version, but we feel that a thorough analysis is beyond
17 the scope of the paper and will be left for future work. 4) *“ReLU, No Biases”*: The choice of Relu with no biases is
18 dictated by requiring interpretability in terms of “input sparsity” of the neural network. In practice, there is no restriction
19 on the network and one can use Relu with biases if input sparsity is not a goal. In any case, removing biases in the
20 ReLU network does not seem to affect performance strongly, especially when the input data is centered.

21 **Reviewer #2.** We thank the reviewer for the positive and encouraging feedback and address their main concerns here.
22 **Additional Comparison on the real word datasets:** For the CCLE experiments, we compared to ElasticNet which is
23 the established feature selection method on this dataset [8, 36]. For the HIV knockoff experiments, we compared to the
24 GLM feature selection of Candès et al [9]. We plan to include the Random Forest feature selections in the final version
25 as an additional baseline, as suggested by the reviewer.

26 **Reviewer #3.** We thank the reviewer for the review and address their questions: 1) *HRT versus Knockoffs*: Thank
27 you for pointing this out and allowing us to clarify that this paper is not a comparative study between the two methods,
28 rather it is a strength of SIC that it can be used with both. For a comparison between HRT and knockoffs, we refer the
29 reviewer to [8], which shows similar performance for either method in terms of controlling FDR. We will highlight in
30 the final paper that each method has its advantages. In HRT most of the computation is in 1) training the generative
31 models, and 2) performing the randomization test, i.e. forwarding the data through the critic and computing p -values
32 for each coordinate for R runs. On the other hand, if knockoff features can be modelled as multivariate Gaussians,
33 controlling FDR with knockoffs can be done very cheaply, since it does not require randomization tests. If instead
34 knockoff features have to be generated through non-linear models, knockoffs can be computationally expensive as well
35 (for example Deep knockoffs, Romano et al. ICLR 2019). 2) *Kernel versus NN SIC*: In early experiments, we tried
36 using random networks defined by random Fourier features (approximating a Gaussian kernel), and we found that a
37 fully trained network outperforms the fixed network. 3) *Computational Complexity versus Performance*: We thank
38 the reviewer for the question, and plan to clarify this point as follows. The cost of training SIC with SGD and mirror
39 descent is of the same order of magnitude as training the base regressor neural network via back-propagation. The only
40 additional overhead is the gradient penalty, where the cost is for a double back-propagation. In our experiments, this
41 added computational cost is not an issue when training is performed on GPU. 4) *Conditions on Φ* : In our framework,
42 Φ is defined on $\mathcal{X} \times \mathcal{Y}$, but is not restricted to be an outer product between feature maps. It is possible to ensure
43 that $SIC(p_{xy}, p_x p_y) = 0$ iff $p_{xy} = p_x p_y$ by imposing universality on the feature map Φ to get injectivity of the
44 mean embeddings. However, studying the conditions on the kernels under which $\eta_j = 0$ leads to $x_j \perp y|x_{-j}$ is
45 an interesting open question, but beyond the scope of this paper. Instead, we tackled the problem (whether η_j are
46 statistically significantly above zero) indirectly using HRT and Knockoffs that have theoretical guarantees to control the
47 FDR with no restriction on the feature selection method used. 5) *Compact level sets of L_ϵ* : This condition is met when
48 perturbing the problem with ϵ . This is not needed for the proof of Theorem 1, but is needed to guarantee that alternating
49 optimization or first order methods on u and η are convergent (See the monograph [19] page 59). We address the minor
50 comments: 1) *Generalized MI*. There are many generalizations of mutual information (MI) such as the Renyi MI that
51 uses Renyi divergence and many other extensions have been developed. For an introduction on this topic we refer
52 the reviewer to α **mutual information** by Sergio Verdu. 2-3) *“zero norm of w , $\frac{1}{N}$ ”* ℓ_0 norm and $\frac{1}{N}$ should be there.
53 Thanks, typos will be corrected. 4) *Multiple runs? GLM?*: This is for a single run, following the same experimental
54 protocol as Candès et al. [9]. GLM stands for Generalized Linear Model. Boosted SIC is an ensemble of η for the
55 same data, different random seeds initialization of the NN. *“L versus g ”* In the cost L the only non-obvious terms for
56 proving convexity in η and u are of the form $g(u, \eta_j) = \frac{u^\top A_j u}{\eta_j}$ (A_j PD) hence the proof is here only for those terms.