

1 We would like to thank the reviewers for their useful, detailed feedback! We will update the paper with the suggested  
2 minor revisions regarding typos and presentation improvements, and respond to individual reviewers comments below.

3 Reviewer #1:

- 4 1. Corollary 1 doesn't actually imply the consistency proven in Theorem 1, because it only applies to a finite  
5 number of moment conditions, whereas in Theorem 1 we deal with an infinite number of moment conditions  
6 and a function space defined by a neural network architecture, which is why proving consistency there is  
7 important and necessary.
- 8 2. You're correct that training the adversarial loss function can be finicky, especially if done naively. We find  
9 that in practice when we include the  $-\frac{1}{4}C$  term inspired by Lemma 1 learning is very stable in particular in  
10 the low dimensional scenarios, and other tweaks like OAdam and early stopping using our validation scheme  
11 seem to improve this further.
- 12 3. We agree that some of the prose discussing GMM+NN results is confusing vis-a-vis results in Table 1. We  
13 will alter this discussion accordingly to more clearly reflect the numbers.
- 14 4. We acknowledge your point about lack of misspecification in DeepIV's first stage in the low-dim scenarios,  
15 and so forbidden regression might not be the right explanation for its performance there. However, the DeepIV  
16 second stage network does have the capacity to fit the problems better than a polynomial. Regarding "full  
17 training": we used an existing implementation of DeepIV.
- 18 5. The network architectures and hyperparams were mistakenly omitted from the supplement and will be added.  
19 Note our code is public. For baseline methods we used the official implementations provided by the authors.
- 20 6. Thanks for catching, we apologize for the typo. Both 0.5's should be 1's:  $X = Z_1 + e + \gamma$ .
- 21 7. Yes, test MSE is w.r.t. the "true response"  $g_0$  over the  $X$  population, as mentioned in 2nd paragraph of Sec 5.

22 Reviewer #3:

- 23 1. Good idea. We will use the notation  $\|x\|_C = x^T C^{-1} x$ .
- 24 2. We will add a proof in the appendix. This is just self-duality of  $L_2$ .
- 25 3. Regarding line 79, our comment is simply pointing out that standard GMM usually starts from the assumption  
26 that the moment conditions specify the problem. Indeed, the assumption that a finite number of moment  
27 conditions identify theta is very strong (too strong) when theta is complex because it easily gives us statistically  
28 efficient methods for estimating theta if true. We will clarify.
- 29 4. Yes; the limit of  $\hat{\theta}_n$  can be anything. This may seem counterintuitive, but one way of understanding this is that  
30 having the "correct" limit corresponds to using the optimally weighted norm, whereas having some other limit  
31 corresponds to some non-optimal norm, but using a non-optimal norm for GMM just means that estimator  
32 won't be statistically efficient, not that it won't be consistent.
- 33 5. All theory in AGMM applies to an *a priori* fixed finite collection of moment conditions; the only "learning [of]  
34 moment functions" in AGMM is in a heuristic jitter step added in the experiments. The significant differences  
35 between DeepGMM and AGMM, other than the drastic difference in performance, are that our method is  
36 inspired by the statistically efficient optimally weighted GMM with a corresponding regularization term which  
37 AGMM lacks, that we directly optimize the neural net critic  $f$ , and that we do learn a critic from an infinite  
38 collection rather than using a finite ensemble of critics. We will update this to make it clearer.
- 39 6. Yes; as line 190 says, what we mean is  $\tilde{\theta}$  is treated as constant. The second term of  $U$  in Equation (9) has zero  
40 partial derivative in  $\theta$ ; so  $\tilde{\theta}$  does not appear in the  $\theta$  gradient. We will rephrase to make this clearer.
- 41 7. This means we use  $f_i(Z) = Z_i$  for  $i = 1, \dots, 784$ . We'll add this clarification.

42 Reviewer #4:

- 43 1. Identification assumption for neural nets: as referenced in line 161, we can easily relax identification and  
44 instead converge to *some*  $\theta$  satisfying the moment conditions (will clarify this simple extension in the proof).  
45 Moreover, this immediately gives that even if there are redundancies in that two  $\theta$ s give the same function  $g$   
46 (e.g., permuting hidden layers), if  $g \in G$  is identified (i.e., all identified thetas give rise to the same function  $g$ ),  
47 we will obtain *some* parameterization  $\theta$  of the unique  $g$ . It's a good point so we will add this discussion.
- 48 2. Regarding the  $-\frac{1}{4}C(f, f)$  term, we found that other regularizers/controls on  $f$  do not perform well as they  
49 induce suboptimal weighting that ignores the covariance of the moment conditions for different  $f$ , whereas  
50 our new regularizer, as you write, is motivated by optimal weighting. This a key driver of our improved  
51 performance over, e.g., AGMM.
- 52 3. The network architectures and learning hyperparams were mistakenly omitted from the supplement and will  
53 be added. Note our code is public. Re "sufficiently rich" moments for Poly2SLS: since its degree is variable,  
54 Poly2SLS can be thought of as a sieve as in Newey and Powell [23], giving universal consistency.
- 55 4. We will cite the recent paper you reference on kernelized 2SLS. It wasn't on our radar (first appearance online  
56 June 2019). We will also cite the Ravuri et al. paper re future work. Thanks.