

1 We thank the reviewers for their thoughtful comments and suggestions and we respond below to some concrete
2 questions/comments that were raised.

3 **Response to Reviewer #4. Slate estimator.** We agree that we should have added a reference to the Swaminathan et
4 al paper. The setup there is also a special case of our setup, where the reward is linear in the treatment vector, i.e.
5 $\langle \theta(z), T \rangle$, but where $T \in R^{\ell \cdot m}$ (ℓ is number of items and m number of slots) and where T takes values in a subset
6 of the hypercube. The discreteness of the action space allows Swaminathan et al to apply a more direct propensity
7 approach (see e.g. Remark 2 for a similar example). Moreover, the slate estimator uses solely the IPS part and is not
8 doubly robust, as the paper works in the setting with a known propensity function.

9 **Efficiency.** This phenomenon is rather standard in the econometrics literature: if one assumes both that the model is
10 well-specified and that there is heteroskedastic noise, then one could typically construct more efficient estimators by
11 optimally re-weighting the samples inversely proportionate to the variance of the error for the corresponding z_i . Such
12 optimally re-weighted estimators are typically avoided in practice as they heavily rely on the well-specification of the
13 model. Hence, we omitted such an analysis. We will add a relevant discussion in the revision with reference of results
14 similar in flavor and an appendix section of how an optimally re-weighted estimator would look like.

15 **Estimating co-variance.** In the worst-case one can view the estimation of the co-variance as a set of separate regressions,
16 one for each entry of the matrix (see e.g. Equation 6 and the sentence above). Assuming the matrix has small dimension
17 and assuming some high-dimensional model space for each regression, then typical regression rates would apply. The
18 example of pricing however shows that in natural problems one might be able to get away with even simpler estimators
19 for this co-variance matrix. We will add some more elaborate discussion on these rates expanding on Remark 4.

20 **Response to Reviewer #5.** See response also to Reviewer #4 regarding efficiency. Roughly: if one assumes the model
21 is well-specified then this implies many more moment conditions than just the unconditional moment implied by a
22 square loss projection. These extra moments can be used to construct more efficient estimators. This can indeed lead
23 to a benefit if the errors are heteroskedastic as then one should do an optimally re-weighted square loss projection.
24 However, these extra moment conditions have no bite in the case of homoskedastic noise. The technical proof in
25 appendix D goes through such an argument. We will add a sketch that highlights these main points.

26 **Response to Reviewer #6. Relation to Foster and Syrgkanis (FS).** Our paper does use the framework of (FS) and the
27 main theorems in that work as a stepping stone in our regret results. However, there are two substantial contributions:
28 1) the framework of (FS) starts from the assumption that one has formulated an orthogonal loss, but does not provide
29 any way of acquiring such an orthogonal loss. Hence, our first contribution is constructing an orthogonal loss in the
30 case of policy learning with continuous actions. Existence of such orthogonal losses were left as an open question in
31 prior work (e.g. Athey and Wager). 2) The out-of-sample regularized ERM provides both a computationally efficient
32 alternative to the variance penalization whenever the original ERM problem is convex and also attains a regret bound
33 whose leading term is much better than that achieved by variance or moment penalization: i) we get a bound that
34 depends on the entropy integral of the policy space as opposed to the critical radius that was achieved by (FS), or the
35 even worse metric entropy at $O(1/n)$ approximation achieved by variance penalization; the latter quantities for instance
36 typically add an extra $\log(n)$ factor in the leading term in the case of VC classes and create even larger deteriorations as
37 compared to the entropy integral for larger classes; achieving an entropy integral dependence has been an open question
38 in the variance penalization literature and the f-divergence robust optimization formulations do not provide an answer
39 to these as they similarly have dependence on metric entropy quantities at fixed approximation levels; moreover the
40 f-divergence equivalence to variance penalization is only asymptotic, ii) we depend on the variance of the difference of
41 the policy loss between the optimal policy and any policy in a small regret slice; this constant can be much smaller than
42 the variance of the optimal policy that is achieved by the moment penalization of (FS) (see discussion after theorem 1).

43 **Computational efficiency.** We note that in both the examples that we present the policy learning problem is convex. We
44 agree that we omitted the description of how we optimize the policy in the pricing example. But we do have a concrete
45 discussion of how we optimize the policy in the costly resource allocation. In the case of pricing, where we optimize
46 over linear policies, then observe that the problem is convex with respect to the coefficients in the linear policy (as it is
47 of the form of maximizing: $\langle \gamma, z \rangle (a(z) + b(z) \langle \gamma, z \rangle)$ and $b(z)$ is non-positive; hence the hessian with respect to gamma
48 is negative semi-definite and hence a concave maximization problem). In this case we optimized the objective by simply
49 finding a closed form solution to the first order condition. This involves simple matrix computations. Similarly as
50 we describe in the costly resource allocation application in Appendix G, the policy learning problem boils down to
51 the square loss minimization over a space of high-dimensional linear policies subject to an ℓ_1 ball constraint (e.g. a
52 multi-task lasso problem); see Equation (55) and statement below. This is a convex problem and can be solved efficiently
53 with standard packages; which is what we employed. In both cases, out-of-sample regularized ERM preserves the
54 convexity of the ERM problem and is efficiently computable via convex optimization; as opposed to variance/moment
55 penalization, which becomes a non-convex problem.