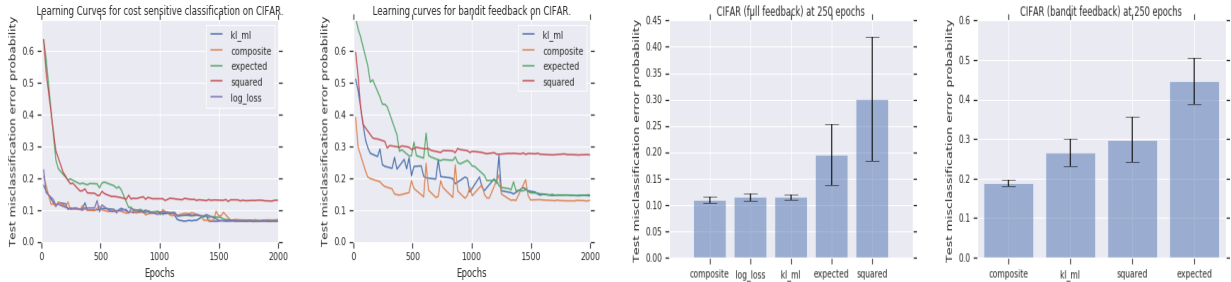


1 **Reviewer #1** Thanks for the comments! $\hat{\mathcal{R}}(\pi)$ was defined in Equation (1), but we will make the reference clearer.

2 We should clarify that the theoretical results already consider out-of-sample generalization. Theorem 8 gives a proper
3 form of generalization bound: with high probability, achieving a small value of the empirical surrogate $\hat{L}(q, \mathcal{D})$
4 guarantees a bound on the **true** suboptimality gap $\mathcal{G}_\tau(\mathbf{f} \circ q)$. The suboptimality gap \mathcal{G}_τ in the smoothed risk \mathcal{S}_τ , defined
5 in (14), is defined in terms of full expectations over the underlying distributions $(x, \mathbf{r}) \sim p(x, \mathbf{r})$ and $a \sim \pi(x)$. In
6 particular, the proof of Theorem 6 explicitly accounts for how the estimator \hat{r} behaves under such a full expectation.

7 There are connections between this work and entropy regularized RL, but there are also distinctions. To be brief: In our
8 scenario, entropy regularization leads to the smoothed risk (5). The connection between smoothed risk and the KL
9 divergence (8) is known. Reversing the KL, as in (9) (without the quadratic), yields a version of maximum entropy
10 inverse RL. A key novelty in this paper is to augment the maxent inverse RL objective with the quadratic, to achieve
11 an *upper bound* on the original KL (8) that remains *calibrated*. In fact, (9) is the key objective we analyze, which is
12 different from (12). This allows us to prove new generalization bounds in the form of Theorem 8. We are also able to
13 achieve successful empirical results for **batch** policy optimization, which remains a challenge in the sequential case
14 (Fujimoto et al. ICML-19). Finally, most work on entropy regularized RL uses split actor-critic models (or considers
15 policy improvement in full MDP planning), whereas we achieve success with a single model that serves as both.

16 **Reviewer #3** Thanks for the comments! We appreciate the criticism of the experimental exposition, and will
17 improve it as suggested, including adding learning curves and standard deviations. Below are the more detailed
18 experimental results (in terms of test misclassification) for CIFAR10, in both the fully observed (§2.2) and partially
19 observed (§3.6) cases. We also have the same suite of results prepared for MNIST, and standard deviations for Table 1.



(a) CIFAR10 fully observed (b) CIFAR10 partially observed (c) CIFAR10 fully observed (d) CIFAR10 partially observed

20 *"It is unclear to me if the reward estimation algorithm is actually evaluated in the experiments."* Yes, Section 3.6 used
21 reward estimation (10), with λ explained in Line 257 and other hyperparameter choices explained in Appendix 4.

22 *"Can you comment on the increased variance demonstrated by Composite on Table 2?"* To produce Table 2,
23 hyperparameters were only chosen to maximize estimated validation reward. We do not yet know whether explicitly
24 considering variance on validation will lead to greater statistical separation between methods.

25 *"I find curious that [...] all the experiments consists of classification tasks "reworked" [...]."* This is inaccurate: the
26 Criteo dataset is a benchmark in this area, which has been extracted from a real online advertising challenge.

27 *"It might also be valuable [...] to run experiments with different data size and distributions other than uniform."* Yes
28 we agree, but to clarify: Table 1 already gives results on CIFAR10 with different training set sizes and non-uniform data
29 collection [12]. The Criteo results in Table 2 are also based on data gathered by a non-uniform logging policy [16].

30 *"Code"* Yes, we plan to release the remaining code. *"Compare to off-policy RL & BayesOpt"* The paper is already
31 performing off-policy RL, albeit for single-step tasks given batch data. The batch scenario implies no exploration, only
32 exploitation. A connection to general RL is given above. We will try to squeeze something in about BayesOpt.

33 **Reviewer #4** Thanks for the comments! The baseline v is chosen as a hyperparameter on validation. It has no
34 effect on the global optima under full expectation, but affects the variance of the empirical estimates.

35 F was actually defined in Line 61, and its use in Proposition 2 was consistent with this definition and subsequent
36 occurrences. We can see how this definition was easily missed however, and we will seek to better highlight it.

37 Thanks for catching the typo in (8). *Line 127:* You are right, the statement should be " $\hat{\mathcal{R}}(\pi)$ yields the highest training
38 error on both MNIST and CIFAR10, the highest test error on CIFAR10, and the second highest test error on MNIST."

39 As far as we are aware, *ImageNet* has never been used as a testbed in this area. It is quite expensive and would only add
40 another "classification" data set to the evaluation. A non-classification based domain might be more interesting.