

1 **Reviewer #4:** Thank you for a detailed review. Regarding your comments on the experiments section, please see the
2 end of this page. We appreciate your comments on the presentation and we will consider incorporating your suggestions
3 in the final version. In terms of content however, we believe that our contributions have been mischaracterized.

4 **Q:** "The contributions of the paper are very close from the one of [12]"

5 **A:** We respectfully strongly disagree. The authors in [12] propose an algorithm following the FTRL scheme, our
6 algorithm is based on minimizing newly proposed surrogate cost functions. They study static regret, we study dynamic
7 regret. They use a generic probabilistic bound at the core of their analysis, we use the specific dynamics of SGD/SGLD.

8 **Q:** "Novelty of the bound on dynamic regret compared to the one of [18]"; "The authors should discuss the links of
9 assumption 3 and the different path-length assumptions in the literature"; "The links with (strong) convexity are not
10 sufficiently discussed"

11 **A:** The regret bound derived in [18] assumes boundedness of the gradients of the cost functions which does not hold for
12 our problem. Furthermore, our assumptions do not imply any of the path-length assumptions in [18]. Our bound and
13 theirs are therefore not comparable. Besides, the work in [18] deals with the general online learning problem, whereas
14 we are only concerned with the specific problem of adaptive importance sampling for SGD/SGLD. There is not much
15 to be said about (strong) convexity apart from mentioning that our cost functions are convex but not strongly convex,
16 which adds to the difficulty of our problem. The whole point of our work is to show that we are able to achieve a good
17 dynamic regret bound even though most of the commonly used assumptions to study dynamic regret do not hold.

18 **Q:** "Why talking about the Langevin method (SGLD) ... it would be just as legitimate as comparing with all stochastic
19 gradient descent methods (SAG, SAGA, SVRG,...)"

20 **A:** It is precisely because the only other method aside from SGD for which our algorithm/analysis work is SGLD that
21 we mention it. In particular, the variance of the SAG/SAGA/SVRG estimators has a completely different form, and
22 designing adaptive importance sampling techniques for these estimators remains an open problem. Furthermore there
23 has been a lot of work in applying variance reduction techniques initially designed for optimization to sampling (see e.g.
24 [35]), and we view our work as a contribution to that literature as well.

25 **Reviewer #6:** Thank you for an in-depth read of our paper. We address your concern below.

26 **Iteration complexity:** At the theoretical level, we can show better convergence rates by replacing our dynamic regret
27 bound directly in the standard convergence bound of averaged SGD as is done in ([11], Theorem 3 and Corollary
28 1). This does not cover the last-iterate case however. Since we are able to obtain per-step regret guarantees (Lemma
29 2 + Lemma 4), we can also obtain bounds on a weighted version of the dynamic regret, which allows us to obtain
30 improved last-iterate guarantees for SGD compared to uniform sampling. We will consider including these results in
31 a new section in the appendix if the reviewer thinks it would be a good addition. In practice, the magnitude of the
32 improvement depends on how much variance reduction can be achieved through the use of the optimal probabilities
33 vs uniform sampling. This is what we addressed in the experiments section, where we mentioned over-parametrized
34 models and imbalanced datasets as examples of cases in which this improvement is significant. **Cost per iteration:**
35 We would like to point out that while our algorithm is $\mathcal{O}(N)$ in both time and space, it has no dimension dependence.
36 For large scale high-dimensional problems, the cost of the final optimization algorithm will still be largely dominated
37 by the gradient evaluations, and the overhead from our sampler will be negligible. The fact that the memory cost is
38 independent of the dimension of the problem is what makes our algorithm an attractive alternative to control-variate
39 methods like SAG/SAGA/SVRG.

40 **Reviewer #7:** Thank you for your very positive comments and helpful suggestions! We will incorporate your
41 suggestions for the experiments section in the final version, please see the paragraph below.

42 **Reviewer #8:** Thank you for taking the time to read our paper.

43 **Concern 1:** Their framework is slightly more general than ours, but when specialized to importance sampling for SGD
44 it becomes exactly the same. Please see the last two paragraphs of the introduction in their paper.

45 **Concerns 2, 3:** See the discussion below. We disagree however on the significance of the theoretical improvement.

46 **Concern 4:** This was a mistake. Please see the supplementary version of the paper.

47 **Comments on the experiments section:** We would like to emphasize that our aim was not to provide a comprehensive
48 empirical study of our algorithm, but only to illustrate its performance on a few datasets, and point to cases in which the
49 improvement is substantial. This is inevitable given the space constraint. That being said, we agree with the reviewers
50 that this section of our work can be improved, and we plan on making full use of the additional page in the final version
51 to do so. In particular we plan on: (i) Giving more details on the setup of the experiments. (ii) Include an experiment
52 that explores more the effect of the choice of step-size (decreasing vs constant) and batch size on the performance of
53 the algorithm. This will include experiments with decreasing step-sizes and a batch-size of 1 which corresponds to the
54 setting of our regret bound. (iii) Include results from other baselines in all of the experiments. (iv) Adding a paragraph
55 analyzing the results and more explicitly explaining when and why the algorithm leads to significant improvements.