

1 Dear reviewers,

2 I would like to thank you all for having taken the time to carefully read through the paper. All three of you were very
3 helpful in pointing out typos and minor errors in the text, as well as providing various suggestions for improvement,
4 which will surely help the paper. I will address each of your comments and concerns individually in separate sections
5 below, although I invite you to read them all.

6 **Reviewer 1**

7 “*The contribution is not clear in the introduction section.*”

8 I agree that the statement of the contributions may not have been as clear as it could have been. I will make the necessary
9 changes so that these contributions are explicitly listed.

10 “*The related work is not sufficiently included in the paper.*”

11 It is indeed possible that the literature survey was not broad enough. If you could provide some additional references, as
12 did reviewer 2, this would be greatly appreciated, both from the point of view of improving the paper and for personal
13 enrichment.

14 “*... How does the rate in (13) extend to the discretized version of the algorithms?*”

15 Although it is possible to compute the bound (13) for each of the models presented in Section 5 in closed form, there is
16 currently no way of *globally* bridging the bound on the continuous dynamics to the discrete case. The attainment of a
17 global bound on these discretized variational methods is the subject of ongoing research. On an *algorithm-by-algorithm*
18 *basis*, bounds on the rates of convergence for almost all of the algorithms presented already have been provided in the
19 works that were cited, so an additional derivation of these were not provided. Based on your concern, which was shared
20 with Reviewer 2, I may include these, and a discussion in relation to (13) in a new version of the paper. I also encourage
21 you to read the last response to Reviewer 3 which may be relevant.

22 **Reviewer 2**

23 “*The stochastic mirror descent problem has also been studied by [1,2] ...*”

24 Thank you very much for these references, they are indeed very relevant. I will include a short discussion of the
25 similarities in a subsequent version of the paper. If you know of any other related works that may have been overlooked,
26 I would greatly appreciate additional references.

27 “*For the discretization methods presented in ... can you guarantee the convergence ...? A discussion ... is needed ...*”

28 Reviewer 1 provided a similar comment. I recommend that you take a look at my response to their third comment
29 above. At a high level, most of the algorithms derived from the models in Section 5 have already been studied in the
30 cited works, and have known rates of convergence associated with them.

31 **Reviewer 3**

32 “*How novel is Theorem 4.1? I am wondering if ... (9)-(11) can be derived using ... techniques ... in Casgrain et al.*”

33 Although the broad ideas behind driving the approach, the techniques in the cited works are specifically developed
34 for linear-quadratic semi-martingale control problems in the context of mean-field games. The set-up in the current
35 paper required some additional theoretical machinery, mainly due to the degree of non-linearity that is present to
36 differentiate through a latent, random Lagrangian. This theorem could not be derived through the direct application of
37 any variational tools were known to me, and I have never before encountered any EL-style equation of this general form
38 in the stochastic control or stochastic variational calculus literature.

39 “*The convergence rate in (13) is ... interesting. Can you comment on ... the stochastic term with time?*”

40 A previous (extended) version of the paper discussed this in more detail. At a high level, we can interpret $d\mathcal{M}_t$ to
41 be the noise introduced through the random sampling of gradients, which causes the stochastic algorithms to deviate
42 from the main effect driven by $d(d\mathcal{L}/d\nu) = (\dots) dt$, corresponding to the deterministic equation of [?]. $\mathbb{E}[[\mathcal{M}]_t]$ can be
43 interpreted as the expected square magnitude of this noise, summed all the way to time t . Thus, the bound (13) tells us
44 that we deviate from the optimal noiseless bound of $O(e^{-\beta t})$ exactly in proportion to how far we expect to deviate
45 from the mean behavior of the algorithm. I hope this answer is what you are looking for, otherwise I would be glad to
46 keep this discussion going.

47 “*The development in Section 5 leaves a lot to be desired because it has a number of debilitating assumptions. ...*”

48 The idea of this section was not to present these models as candidates for realistic models of the loss function and
49 observation dynamics. Rather, the point was to answer the question: *For an existing discrete stochastic optimization*
50 *algorithm (e.g. SGD or stochastic momentum), what are the implicit assumptions made by this algorithm on the*
51 *optimization problem it is trying to solve, and under what problem conditions is the algorithm ‘optimal’ in the sense*
52 *of the variational model?* It is indeed a surprising result that very commonly used stochastic optimization algorithms
53 implicitly make these very simplistic assumptions about the problem structure.

54 I believe that the motivation for this section can be made more clear, since it may confuse readers in the current state.