

1 We thank all the reviewers for the helpful comments and questions. Before moving on to answer reviewers’ questions,  
2 we first clarify the main **contributions** and **motivations** of our paper.

3 The goal of the paper is to better understand neural network optimization, especially the interactions between algorithmic  
4 choices and batch-size. We proposed a simple, yet useful toy model – Noisy Quadratic Models. This model captures  
5 the essential behavior we see in real neural networks (as validated by our experiments) while **allowing us to run**  
6 **experiments in seconds**. Note that most assumptions we made are based on existing theoretical results or experiments  
7 in real neural networks. We stress that NQM makes it easy to test new empirical ideas for practitioners and derive new,  
8 testable theoretical results for theorists (as an example, we studied the role of momentum analytically, showing that the  
9 momentum and learning rate are interchangeable in the small batch regime, which is non-trivial in stochastic setting).

10 We now address all comments and questions in order.

11 **Interaction of momentum and preconditioning (R1).** For neural net experiments, we applied preconditioning before  
12 the momentum. For the NQM, with fixed preconditioning, both methods are equivalent. We will update the paper to  
13 clarify these points.

14 **Relationship between the Hessian and Fisher matrices (R1).** For the models we study, the Fisher matrix is equivalent  
15 to the Gauss-Newton Hessian; see Martens, (2014) for details. We will make this explicit in the paper.

16 **Assumption of diagonal Hessian (R2, R3).** The assumption that the Hessian is diagonal can be made *without loss of*  
17 *generality* in the sense that (as noted in lines 94–98) we focus on algorithms such as SGD and Heavy-ball momentum  
18 that are invariant under variable rotations. For such algorithms, we can analyze the evolution of iterates in a basis that  
19 makes the Hessian diagonal, without changing the dynamics of the system.

20 **Theoretical contribution (R2).** We stress here that the analysis of momentum in stochastic setting is non-trivial. In  
21 particular, stochasticity (gradient noise) introduces extra difficulty in analyzing heavy-ball momentum. For the NQM,  
22 we showed theoretically that momentum SGD performs similarly to plain SGD in the regime of small batch sizes but  
23 helps in the large-batch regime, which also matches previous studies and our large-scale experiments.

24 **NQM is too simple, and some assumptions are false for neural nets (R2, R3).** When constructing a model for some  
25 phenomenon, the model does not need to be completely faithful in all respects; rather, it is the act of abstracting away  
26 inessential features that allows one to tractably analyze phenomena. Simplicity is a virtue, as long as the model captures  
27 the effects relevant to the phenomenon. One of our main contributions is showing that (empirically) features such as  
28 non-convexity, non-stationary gradient noise, etc., are not needed to explain certain neural net training phenomena  
29 which have recently received a lot of attention. Abstracting away these features will make it much easier for others  
30 to build on our work by further analyzing these phenomena. (Of course, for explaining other phenomena we don’t  
31 consider, such as local optima, one might need to re-introduce features such as non-convexity.)

32 **Assumptions of fixed gradient noise and particular form of preconditioner (R2).** (Note that we *don’t* require the  
33 noise distribution to be Gaussian.) Both assumptions are required for analytic analysis and fast simulation, yet reflect  
34 (to some extent) the reality of neural networks. Particularly, we showed in the appendix that the Fisher matrix stabilizes  
35 after a few epochs of training on ResNet, supporting the assumption of fixed gradient noise. In other words, we expect  
36 the assumption fixed gradient noise roughly captures the essence of real neural networks in throughout all of training  
37 except the initial phase.

38 **Quadratic loss function (R2, R3).** We stress here that there are several ways one might try to justify the use of  
39 convex quadratic models of the objective (also see lines 83-89 of our submission). First is the well-known fact that  
40 any smooth function will resemble a convex quadratic in a small enough neighborhood around a local minimizer. And  
41 recent theoretical work (Jacot et al., 2018, Du et al., 2018, Lee et al., 2019) has argued that very wide neural networks  
42 will deviate only a small distance from the initial random point in parameter space throughout training, and that they  
43 thus behave similarly to linearized networks for the purposes of training. This implies that the objective will appear  
44 “locally convex” throughout training, and that for the squared error loss will even resemble a convex quadratic. Adding  
45 further support for the “networks as linearized models” approximation is the fact that the Generalized Gauss-Newton  
46 matrix (Martens, 2014), which is the matrix of choice behind the most powerful 2nd-order neural network optimization  
47 methods, can be seen as the Hessian of the training objective under said approximation. Finally, we note that quadratic  
48 approximations also motivated Laplace approximation (MacKay, 1992) and variational inference (Zhang et al., 2018) in  
49 Bayesian neural networks.

50 We thank all reviewers again. We hope that our responses address your comments and concerns.