

1 We thank all three reviewers for their thorough reviews and constructive feedback.

2

3 **Why not Newton?:**

4 As remarked by reviewers #6 and #7, retaining the diagonal blocks of the Hessian in the approximation results in a
5 regularized and damped Newton method $(\Delta x, \Delta y) = (\text{Id} + \eta J)^{-1}(\nabla_x f, \nabla_y g)$, where J is the Jacobian of the vector
6 field $(\nabla_x f, \nabla_y g)$. We agree that the present version of the paper lacks a thorough discussion of the reasons for ignoring
7 the diagonal parts of the Hessian and will therefore add the following three reasons.

8 **Blow-up of condition number:** Including the diagonal blocks of the Hessian in the non-convex-concave setting can
9 make the matrix inverse arbitrarily ill-conditioned as $\|\eta D_{xx}^2 f\|$ or $\|\eta D_{yy}^2 f\|$ approach or exceed 1, greatly increasing
10 the cost of the linear system solve. In contrast, for zero-sum games the condition number of the matrix inverse in CGD
11 is always bounded from above by $(1 + \eta^2 \|D_{xy}^2 f\|^2)$.

12 **Irrational updates:** For $\eta \|D_{xx}^2 f\|$ or $\eta \|D_{yy}^2 f\|$ bigger than 1 and f non-convex-concave, the regularized Newton
13 update can lose its game-theoretic interpretation as a local strategic equilibrium, allowing for convergence to highly
14 non-optimal critical points. While we leave the full characterization of the attractors of CGD for future work, we expect
15 them to always be game-theoretically meaningful since the updates of CGD arise as local Nash equilibria.

16 **Lack of regularity:** For the diagonal blocks of the Hessian to be useful in optimization, we need to make additional
17 assumptions on the regularity of the loss function, for example by bounding the Lipschitz constants of $D_{xx}^2 f, D_{yy}^2 f$.
18 Otherwise, including additional second order information can make the results worse. Consider for instance, minimizing
19 $x \mapsto x^2 + \epsilon^{3/2} \sin(x/\epsilon)$ for $\epsilon \ll 1$. Many minimax problems, for example GANs, have the form $f(x, y) =$
20 $\Phi(\mathcal{G}(x), \mathcal{D}(y))$ where Φ is *smooth* and *simple* but \mathcal{G} and \mathcal{D} might only have first order regularity. In this setting, the
21 bilinear approximation has the advantage of fully exploiting first order information of \mathcal{G}, \mathcal{D} , without assuming them to
22 have higher degrees of regularity. This is because the bilinear approximation of f then contains only the first derivatives
23 of \mathcal{G} and \mathcal{D} , while the quadratic approximation contains second derivatives $D_{xx}^2 \mathcal{G}$ and $D_{yy}^2 \mathcal{D}$, and therefore needs
24 stronger regularity assumptions on \mathcal{G} and \mathcal{D} to be effective.

25 **Reviewer #3:**

26 **"...convergence rate results for CGD...":** Under lower (upper) bounds on $D_{xx} f (D_{yy} f)$, global exponential conver-
27 gence can be derived from Theorem 2.2, and we are happy to include this result with the revisions. A special case of this
28 is strong convex-concavity. We are not aware of existing work on minimaximization that provides global convergence
29 proofs without either convex-concavity/monotonicity, or strong additional assumptions.

30 **"...CGD still requires that the step-size is bounded by one over the max diagonal entry of the Hessian...":** Correct!
31 This is the analogous requirement to applying gradient descent to the single player game (keeping the other player
32 fixed). For problems like GANs the problem of optimizing one player while keeping the other player fixed can be
33 solved reliably via gradient descent while the two-player game becomes unstable under alternating gradient descent.
34 The purpose of CGD is to solve two-player games with similar step sizes and stability properties as when using gradient
35 descent to optimize one player while the other player is kept fix.

36 **Reviewer #6:**

37 **Concern 1: Why not use full second order?** See first paragraph.

38 **Concern 2: Complexity of matrix inverse?** We are sorry for the misunderstanding and will try to make line 289 more
39 precise: Figure 5 does **not** show the convergence as a function of the iteration count, but as a function of the number
40 of gradient evaluations and Hessian-vector products. Thus, a single step of CGD that needs k iterations of conjugate
41 gradient to solve the linear system in the update rule will amount to an x -value of $(4 + 2k)$, while a single step of
42 optimistic gradient descent ascent (OGDA) corresponds to an x -value of 2 in the plot. Thus, this measure of cost fairly
43 accounts for the complexity of the matrix inverse in CGD. We find that the convergence rate of CGD is competitive
44 throughout a wide range of step sizes and we explain this fact with the *graceful reduction to linearized CGD* described
45 in Line 274: If the gain from the matrix inverse is small, the matrix will be well-conditioned and thus easy to invert.
46 See also our answer to Reviewer #7.

47 **Reviewer #7:**

48 **Concerns 1 & 2: Why drop diagonal blocks of Hessian? Why use bilinear approximation?** See first paragraph.

49 **Concern 3: Is CGD scalable?** Since mixed mode automatic differentiation allows to compute Hessian vector products
50 with minimal overhead compared to gradient computations using reverse mode automatic differentiation, we see no
51 reason why CGD should be restricted to small problems. While larger problems will tend to require more iterations of
52 conjugate gradient, they also tend to negatively affect the existing methods, with the experiments in Figure 5 suggesting
53 that the advantage of CGD over existing methods *increases* as the problems size increases. We are presently working
54 on an implementation of CGD using JAX (which provides GPU accelerated mixed mode automatic differentiation) to
55 then apply to large GAN problems.