

1 We thank reviewers for the time they took to read the paper and provide comments and suggestions.

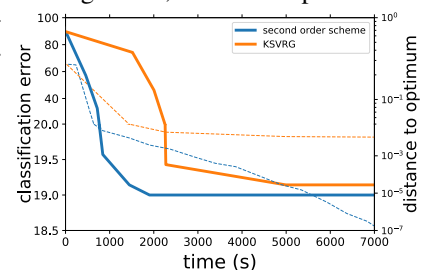
2 **Reviewer 1.** As the reviewer correctly points out one crucial aspect of the proposed algorithm is exactly the inde-
3 pendence from the condition number of the problem. More generally, note that the paper provides an optimization
4 algorithm that at the same time is: (a) independent (except for logarithmic terms) from the condition number of the
5 problem, (b) globally convergent (unlike [6,7], the method will reach the optimum point from zero), (c) efficient from a
6 computational viewpoint. Essentially all the known optimization algorithms nowadays satisfy only a subset of the points
7 above. Indeed first order methods don't satisfy (a); Second order methods for general losses only (a); Second order with
8 approximation methods as in Def. 2 satisfy (a)+(c). However the fact that they are not globally convergent make them
9 not usable on real problems, since a first order method is needed to arrive very close to the optimum, losing (a). Finally
10 there exists the very peculiar class of self-concordant function for which it has been proven that it is possible to achieve
11 (a)+(b)+(c). However such result is more of theoretical interest in ML since very few objective functions there are self
12 concordant. The novel contribution of the paper is a globally convergent optimization scheme that allows to achieve
13 (a)+(b)+(c) for a larger class of functions i.e. *generalized self concordant* functions (GSC), which instead contains many
14 losses of interest for regression, multivalued regression, robust regression, multiclass classification and multilabeling.

15 *Relevance of GSC for ML and kernel methods.* We would like to stress the significance of our result for part of ML
16 and related fields. Indeed, while GSC is a subset of convex functions, it still covers many widely used approaches in
17 machine learning, model estimation, statistics and finance, which are still framed in terms of large scale generalized
18 linear models, like logistic and soft-max regression. More specifically, kernel regression and classification methods
19 have the peculiar property to lead to large scale convex problem even if they are learning non-linear functions. In
20 the community of kernel methods it is still an open problem how to find efficient methods to solve such big convex
21 problems, without depending on the condition number, and our result is can constitute a key step in this direction.

22 *Second order skepticism.* We also understand the skepticism towards second order methods. Note however that (1) they
23 become crucial to solve severely ill-conditioned problems, where the condition number is $\gg n$ i.e. 10^{10} , since first
24 order methods can't avoid the dependence on the condition number. (2) the fact that second order methods have been
25 labeled as inefficient in the '90 exactly because proper approximation techniques *and* globalization schemes were
26 missing. In the past few years many approximation schemes have been developed as discussed in the introduction using
27 novel results from randomized linear algebra. Nowadays as discussed in Section 2, the cost per iteration (epoch) of
28 one approximated second order method is essentially the same as gradient descent (or stochastic gradient descent).
29 Crucially the globalization scheme provided in this paper guarantees that the number of required iterations does not
30 depend polynomially on the condition number (as in gradient descent or the stochastic methods where it is linear, or a
31 square root with acceleration), but only logarithmically. So we need way less iterations than a first order method.

32 In particular we will make the introduction more clear putting elements of the discussion above; we will add a paragraph
33 after the definition of GSC function to clarify their role; we will summarize the computational costs of the various
34 algorithms of Section 2 via an additional table and moved some of the crucial aspects on how we improved over the
35 state of the art from Section 2 to the contribution paragraph in the spirit of the discussion above.

36 *Experiments.* The proposed paper is clearly of theoretical flavour. We provide a new algorithm, and crucial part of the
37 contribution is the development of the mathematical theory in terms of techniques, theorems and proofs to guarantee the correct behaviour of the algorithm.
38 Additionally to show the practical validity of the approach, we performed a large
39 scale experiment in the field of kernel methods (see above). In particular, in the
40 experiments we compared our approach with KSVRG that is nowadays one of
41 the fastest accelerated first order stochastic methods and is widely used to train
42 kernel methods. We will add this plot showing the classification error for the
43 experiment, which shows interestingly that it is still hard for first order methods,
44 to achieve good performance in classification, for ill conditioned problems.



46 **Reviewer 2.** We thank Reviewer 2 the very thorough reading of the paper and for pointing out many typos which we
47 will correct. We agree that FALKON is a very related algorithm; Indeed, the aim of the paper is to extend FALKON
48 to more general losses. In particular the proposed approach recovers it in the case of squared loss. We are currently
49 working on a minimal usable implementation for python and multicore-GPU, we will make the code open source and
50 available on github this fall. Finally we will correct the format of the equation and the typos as suggested.

51 **Reviewer 3.** We thank Reviewer 3 for the very thorough understanding of the paper. We will add a note about
52 the difference between GSC [6] and self concordant functions [24] after the definition of GSC functions. The
53 difference with Newton sketch and our paper is that Newton sketch's globally convergent scheme applies to self-
54 concordant functions while ours applies to GSC functions. Moreover as suggested, in Section 1.1 we will emphasize
55 the advantages/disadvantages of our method. Indeed, its aim is to extend the performance of quadratic solvers to GSC
56 functions, up to a log factor (only quadratic solvers are "condition number independent").