

1 We would like to thank the reviewers for their detailed and insightful comments.

2 **[Performance of SGD and robustness of optimization algorithms.]** We have resolved the concerns with SGD. By  
 3 increasing the batch size towards the last iterations and averaging the last iterates, SGD on the adaptive Gaussian sketch  
 4 problem performs better in terms of time vs accuracy performance, compared to SGD on problem (1) or on the oblivious  
 5 Gaussian sketch problem. We have similar results with Adam. As reported in the submission, SVRG on the adaptive  
 6 Gaussian sketch performs better than SVRG on problem (1), and is robust to the choice of hyperparameters. Further,  
 7 Sub-sampled Newton (with mini-batch Hessian and full-batch gradient) has a strong time vs accuracy performance on  
 8 adaptive Gaussian sketch. In the revised version, we will include our new results for SGD and Adam, and a sensitivity  
 9 analysis to sketching, batch and step sizes, for all algorithms applied to the sketched problems (adaptive and oblivious).

10 **[Comparison with other sketching baselines.]** We carried out extensive numerical evaluations of oblivious Gaussian  
 11 sketching and adaptive sketching with uniform column sub-sampling matrix (Nystrom method) on MNIST and CIFAR10.  
 12 For a wide range of values of sketching size  $m$  and regularization parameter  $\lambda$ , adaptive Gaussian sketching always  
 13 strongly beats oblivious sketching, and, outperforms Nystrom method, both in terms of final test accuracy (see Table 1  
 14 below), and, time vs accuracy performance for the following algorithms: SGD, SVRG, Sub-sampled Newton and Adam.  
 15 Further, adaptive Gaussian sketching matches the performance of  $x^*$  for relatively small values of  $m$ . We will include  
 all these results in the revised version. **[Computational issues with  $(S^\top S)^{-\frac{1}{2}}$  for large  $m$ .]** Thanks to this question,

Table 1: Test classification error on MNIST and CIFAR10, for 10-classes classification. "AG": Adaptive Gaussian  
 sketch, "Ob": Oblivious Gaussian sketch, "N": Nystrom method,  $x_m$ : solution obtained from problem (2) with sketching  
 size  $m$ . We mapped MNIST (resp. CIFAR10) images through 10000 (resp. 60000) random cosines.

$\lambda$	$x_{MNIST}^*$	$x_{256}^{AG}$	$x_{1024}^{AG}$	$x_{256}^{Ob}$	$x_{1024}^{Ob}$	$x_{256}^N$	$x_{1024}^N$	$x_{CIFAR}^*$	$x_{256}^{AG}$	$x_{1024}^{AG}$	$x_{256}^{Ob}$	$x_{1024}^N$	$x_{256}^N$	$x_{1024}^N$
$5 \cdot 10^{-5}$	4.6 %	4.0 %	4.5 %	25.2 %	8.5%	5.0 %	4.6 %	51.6 %	50.6%	51.0%	70.5%	55.8%	53.1%	
$5 \cdot 10^{-6}$	2.5%	2.8%	2.4%	30.1%	9.4%	3.0%	2.7%	47.6%	51.9%	45.8%	80.1%	57.2%	55.8%	

16  
 17 we have improved our results and we can show that the matrix  $(S^\top S)^{-\frac{1}{2}}$  can be replaced by any other pre-conditioner  
 18  $Q$ , and in particular, for large  $m$ , a matrix  $Q$  obtained by approximate SVD. Provided  $\|Q - (S^\top S)^{-\frac{1}{2}}\|_2$  is small, then  
 19 the condition number of (7) remains close to that of (1). Importantly, it does not affect any of the bounds on  $\tilde{x}$ . We will  
 20 include these results in the revised version.

21 **[For small  $m$ , would dynamically modifying the sketching matrix lead to tighter bounds?]** For small  $m$ , we tried  
 22 numerically to refresh the sketching matrix at each iteration and it did not yield good results. However, our Algorithm 2  
 23 refreshes the sketching matrix at the end of each optimization, and gives tighter bounds.

24 **[Results on CIFAR10 far from state-of-the-art. Other optimization problems for which the method could be  
 25 demonstrated?]** We did additional experiments with features extracted from a pre-trained neural network, and  $\tilde{x}$   
 26 matched exactly the test error of  $x^*$  ( $\sim 10\%$ ). We will include these results in the final version. Beyond classification,  
 27 large-scale generalized linear models (other than least squares) can be addressed with our method.

28 **[Unclear if SGD (without sketching) converges to the same solution or performs better].** The reported results  
 29 correspond to the best SGD solution we obtained (with grid search of the batch and step sizes), even through longer  
 30 time horizons.

31 **[Value of  $\lambda$  used for synthetic experiments?]** We used  $\lambda = 10^{-4}$ . Thank you for pointing this out, we will fix it.

32 **[Title suggestion.]** We agree with the relevant title suggestion. **[Non-convexity.]** We derived a new result regarding  
 33 non-convex, smooth functions  $f$ : if  $\alpha^*$  is a nearly-stationary point for the sketched problem (2), then  $\tilde{x}$  is a nearly  
 34  $\varepsilon$ -stationary point for problem (1), where  $\varepsilon$  controlled again by the tail spectral decay of  $A$ . **[Regularity assumptions  
 35 on  $f$ .]** We will add some discussion in the main body of the paper. In Appendix E, guarantees are provided for convex,  
 36 non-smooth objectives  $f$ . **[Other regularizers.]** We have extended our analysis to smooth, strongly convex regularizers.  
 37 However, extension to the  $L_1$ -norm is an open question.

38 **[Invoking the representer theorem not necessary in Eq. (11).][Notation 5.10 $^{-5}$  non-standard.]** We will simplify  
 39 the argument for Eq. (11) in the revised version, and correct the notation. Thank you for pointing this out.

40 **[Intuition for why the proposed approach works.]** We will discuss more carefully some intuition in the revised  
 41 version. In a nutshell, the kernelized version of optimization problem (1) is well approximated by  $\min_w f(AP_S A^\top w) +$   
 42  $\lambda \|P_S A^\top w\|^2$ , provided that  $AA^\top \approx AP_S A^\top$ . Adaptive sketching works better than oblivious one, since  $\|AA^\top -$   
 43  $AP_{A^\top \tilde{S}} A^\top\|_2 \ll \|AA^\top - AP_{\tilde{S}} A^\top\|_2$ , for  $\tilde{S}$  i.i.d. Gaussian.

44 **[Comparison of the proposed method with approximate kernel methods]** Random Fourier features lead to problems  
 45 of type (1). Standard Nystrom methods approximates the matrix  $K$  in (11). But both problems (1) and (11) are typically  
 46 high-dimensional. Our method is a dimension-reduction tool, that can be used on top of approximate kernel methods.