
Supplement: Competitive Gradient Descent

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This is the supplement to the paper "Competitive Gradient Descent"

2 1 Proofs of convergence

3 *Proof of Theorem 2.3.* To shorten the expressions below, we set $a := \nabla_x f(x_k)$, $b := \nabla_y f(x_k, y_k)$,
 4 $H_{xx} := D_{xx}^2 f(x_k, y_k)$, $H_{yy} := D_{yy}^2 f(x_k, y_k)$, $N := D_{xy}^2 f(x_k, y_k)$, $\tilde{N} := \eta N$, $\tilde{M} := \tilde{N}^\top \tilde{N}$, and
 5 $\bar{M} := \tilde{N} \tilde{N}^\top$. Letting (x, y) be the update step of CGD and using Taylor expansion, we obtain

$$\begin{aligned} & (\nabla_x f(x + x_k, y + y_k))^2 + (\nabla_y f(x + x_k, y + y_k))^2 - \|a\|^2 - \|b\|^2 \\ & \leq 2x^\top H_{xx} a + 2x^\top N b + 2a^\top N y + 2b^\top H_{yy} y \\ & \quad + 4L(\|x\|^2 + \|y\|^2)(\|a\| + \|b\|) \\ & = +2\eta \left(-a^\top - b^\top \tilde{N}^\top \right) (\text{Id} + \bar{M})^{-1} H_{xx} a \\ & \quad + 2x^\top N b + 2a^\top N y \\ & \quad + 2\eta b^\top H_{yy} \left(\text{Id} + \tilde{M} \right)^{-1} \left(b - \tilde{N}^\top a \right) \\ & \quad + 4L(\|x\|^2 + \|y\|^2)(\|a\| + \|b\|) = \dots, \end{aligned}$$

6 By expanding zero to $\pm 2\eta b^\top \tilde{N}^\top (\text{Id} + \bar{M})^{-1} H_{xx} a$ and $\pm 2\eta b^\top H_{yy} \left(\text{Id} + \tilde{M} \right)^{-1} \tilde{N}^\top a$, we obtain

$$\begin{aligned} \dots & = -2\eta a^\top H_{xx} a + 2\eta a^\top \bar{M} (\text{Id} + \bar{M})^{-1} H_{xx} a \\ & \quad - 2\eta b^\top \tilde{N}^\top (\text{Id} + \bar{M})^{-1} H_{xx} a \\ & \quad + 2x^\top N b + 2a^\top N y \\ & \quad + 2\eta b^\top H_{yy} b + b^\top H_{yy} \left(\text{Id} + \tilde{M} \right)^{-1} \tilde{M} b \\ & \quad - 2\eta b^\top H_{yy} \left(\text{Id} + \tilde{M} \right)^{-1} \tilde{N}^\top a \\ & \quad + 4L(\|x\|^2 + \|y\|^2)(\|a\| + \|b\|) = \dots \end{aligned}$$

7 We now plug the update rule of CGD into x and y and observe that $\tilde{N}^\top (\text{Id} + \bar{M})^{-1} =$
 8 $(\text{Id} + \tilde{M})^{-1} \tilde{N}^\top$ to obtain

$$2x^\top N b + 2a^\top N y = -2a^\top (\text{Id} + \bar{M})^{-1} \bar{M} a - 2b^\top \left(\text{Id} + \tilde{M} \right)^{-1} \tilde{M} b.$$

9 By plugging this into our main computation, we obtain

$$\begin{aligned}
\dots &= -2\eta a^\top H_{xx}a + 2\eta a^\top \bar{M} (\text{Id} + \bar{M})^{-1} H_{xx}a \\
&\quad - 2\eta b^\top \tilde{N}^\top (\text{Id} + \bar{M})^{-1} H_{xx}a \\
&\quad - 2a^\top (\text{Id} + \bar{M})^{-1} \bar{M}a - 2b^\top (\text{Id} + \tilde{M})^{-1} \tilde{M}b \\
&\quad + 2\eta b^\top H_{yy}b - 2\eta b^\top H_{yy} (\text{Id} + \tilde{M})^{-1} \tilde{M}b \\
&\quad - 2\eta b^\top H_{yy} (\text{Id} + \tilde{M})^{-1} \tilde{N}^\top a \\
&\quad + 4L(\|x\|^2 + \|y\|^2)(\|a\| + \|b\|) \leq \dots
\end{aligned}$$

10 By positivity of squares, we have

$$\begin{aligned}
2\eta a^\top \bar{M} (\text{Id} + \bar{M})^{-1} H_{xx}a &\leq a^\top \left(\bar{M} (\text{Id} + \bar{M})^{-1} \right)^2 a + a^\top (\eta H_{xx})^2 a \\
-2\eta b^\top H_{yy} (\text{Id} + \tilde{M})^{-1} \tilde{M}b &\leq b^\top \left(\tilde{M} (\text{Id} + \tilde{M})^{-1} \right)^2 b + b^\top (\eta H_{yy})^2 b.
\end{aligned}$$

11 For $\lambda \in [-1, 1]$ we have $-2\lambda + \lambda^2 = 2\lambda(1 - \lambda/2) \leq -h_\pm(\lambda)$ from which we deduce the
12 result. \square

13 Theorem 2.4 follows from Theorem 2.3 by relatively standard arguments:

14 *Proof of Theorem 2.4.* Since $\nabla_x f(x^*, y^*), \nabla_y f(x^*, y^*) = 0$ and the gradient and Hessian of f are
15 continuous, there exists a neighbourhood \mathcal{V} of (x^*, y^*) such that for all possible starting points
16 $(x_1, y_1) \in \mathcal{V}$, we have $\|(\nabla_x f(x_2, y_2), \nabla_y f(x_2, y_2))\| \leq (1 - \lambda_{\min}/4)\|(\nabla_x f(x_1, y_1), \nabla_y f(x_1, y_1))\|$.
17 Then, by convergence of the geometric series there exists a closed neighbourhood $\mathcal{U} \subset \mathcal{V}$ of (x^*, y^*) ,
18 such that for $(x_0, y_0) \in \mathcal{U}$ we have $(x_k, y_k) \in \mathcal{V}, \forall k \in \mathbb{N}$ and thus (x_k, y_k) converges at an
19 exponential rate to a point in \mathcal{U} . \square

20 2 Details regarding the experiments

21 2.1 Experiment: Estimating a covariance matrix

22 We consider the problem $-g(V, W) = f(W, V) = \sum_{ijk} W_{ij} (\hat{\Sigma}_{ij} - (V\hat{\Sigma}V^\top)_{i,j})$, where the $\hat{\Sigma}$
23 are empirical covariance matrices obtained from samples distributed according to $\mathcal{N}(0, \Sigma)$. For our
24 experiments, the matrix Σ is created as $\Sigma = UU^\top$, where the entries of $U \in \mathbb{R}^{d \times d}$ are distributed
25 i.i.d. standard Gaussian. We consider the algorithms OGD, SGA, ConOpt, and CGD, with
26 $\gamma = 1.0, \epsilon = 10^{-6}$ and let the stepsizes range over $\eta \in \{0.005, 0.025, 0.1, 0.4\}$. We begin with the
27 deterministic case $\hat{\Sigma} = \Sigma$, corresponding to the limit of large sample size. We let $d \in \{20, 40, 60\}$
28 and evaluate the algorithms according to the trade-off between the number of forward evaluations and
29 the corresponding reduction of the residual $\|W + W^\top\|_{\text{FRO}}/2 + \|UU^\top - VV^\top\|_{\text{FRO}}$, starting with
30 a random initial guess (the same for all algorithms) obtained as $W_1 = \delta W, V_1 = U + \delta V$, where the
31 entries of $\delta W, \delta V$ are i.i.d uniformly distributed in $[-0.5, 0.5]$. We count the number of "forward
32 passes" per outer iteration as follows.

- 33 • OGD: 2
- 34 • SGA: 4
- 35 • ConOpt: 6
- 36 • CGD: $4 + 2 \times \text{number of CG iterations}$

37 The results are summarized in Figure 1. We see consistently that for the same stepsize, CGD has
38 convergence rate comparable to that of OGD. However, as we increase the stepsize the other
39 methods start diverging, thus allowing CGD to achieve significantly better convergence rates by
40 using larger stepsizes. For larger dimensions ($d \in \{40, 60\}$) OGD, SGA, and ConOpt become

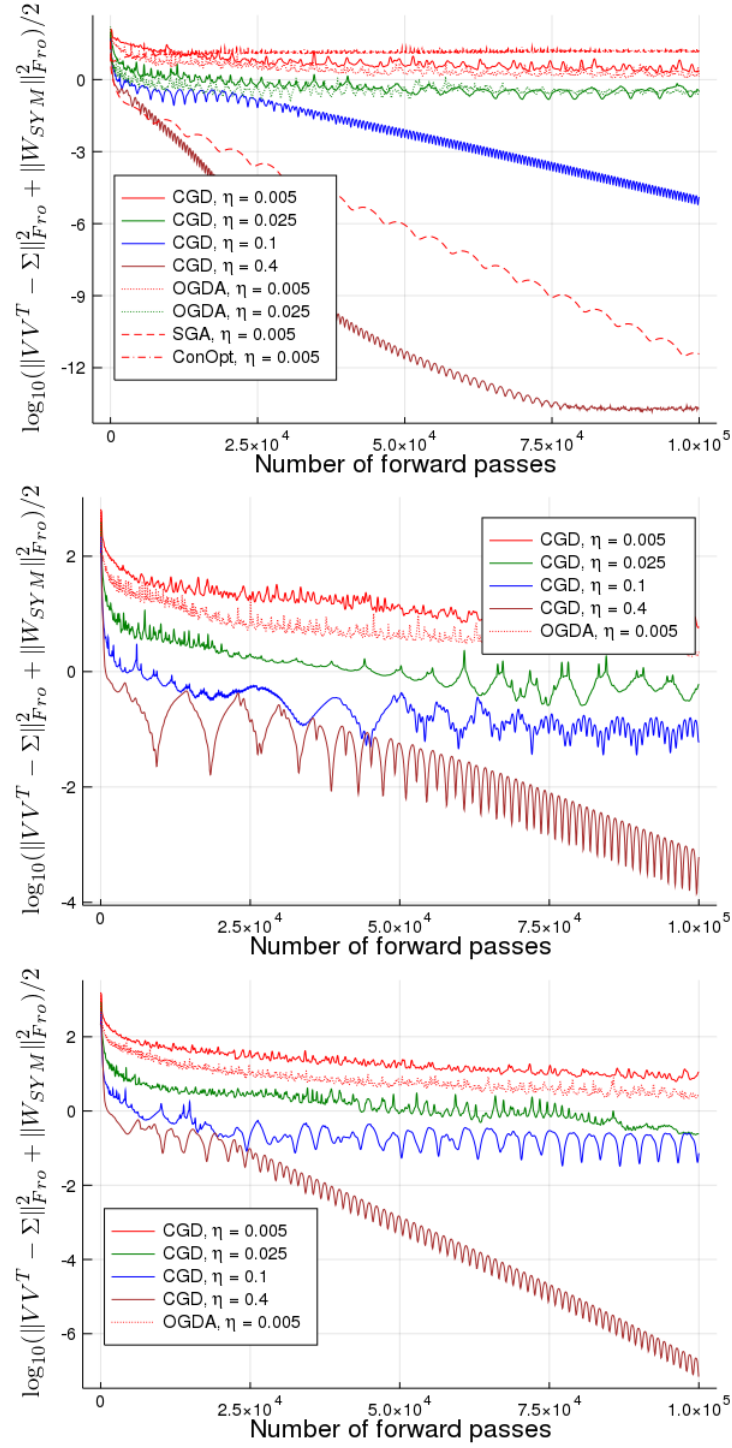


Figure 1: The decay of the residual as a function of the number of forward iterations ($d = 20, 40, 60$, from top to bottom). **Note that missing combinations of algorithms and stepsizes correspond to divergent experiments.** While the exact behavior of the different methods is subject to some stochasticity, results as above were typical during our experiments.

41 even more unstable such that OGDA with the smallest stepsize is the only other method that still
 42 converges, although at a much slower rate than CGD with larger stepsizes. We now consider the
 43 stochastic setting, where at each iteration a new $\hat{\Sigma}$ is obtained as the empirical covariance matrix
 44 of N samples of $\mathcal{N}(0, \Sigma)$, for $N \in \{100, 1000, 10000\}$. In this setting, the stochastic noise very
 45 quickly dominates the error, preventing CGD from achieving significantly better approximations than
 46 the other algorithms, while other algorithms decrease the error more rapidly, initially. It might be
 47 possible to improve the performance of our algorithm by lowering the accuracy of the inner linear
 48 system solve, following the intuition that in a noisy environment, a very accurate solve is not worth
 49 the cost. However, even without tweaking ϵ it is noticeable that the trajectories of CGD are less noisy
 50 than those of the other algorithms, and it is furthermore the only algorithm that does not diverge for
 51 any of the stepsizes. It is interesting to note that the trajectories of CGD are consistently more regular
 52 than those of the other algorithms, for comparable stepsizes.

53 **2.2 Experiment: Fitting a bimodal distribution**

54 We use a GAN to fit a Gaussian mixture of two Gaussian random variables with means $\mu_1 = (0, 1)^\top$
 55 and $\mu_2 = (2^{-1/2}, 2^{-1/2})^\top$, and standard deviation $\sigma = 0.1$. Generator and discriminator are given
 56 by dense neural nets with four hidden layers of 128 units each that are initialized as orthonormal
 57 matrices, and ReLU as nonlinearities after each hidden layer. The generator uses 512-variate standard
 58 Gaussian noise as input, and both networks use a linear projection as their final layer. At each
 59 step, the discriminator is shown 256 real, and 256 fake examples. We interpret the output of the
 60 discriminator as a logit and use sigmoidal crossentropy as a loss function. We tried stepsizes
 61 $\eta \in \{0.4, 0.1, 0.025, 0.005\}$ together with RMSProp ($\rho = 0.9$) and applied SGA, ConOpt ($\gamma = 1.0$),
 62 OGDA, and CGD. Note that the RMSProp version of CGD with diagonal scaling given by the
 63 matrices S_x, S_y is obtained by replacing the quadratic penalties $x^\top x/(2\eta)$ and $y^\top y/(2\eta)$ in the
 64 local game by $x^\top S_x^{-1}x/(2\eta)$ and $y^\top S_y^{-1}y/(2\eta)$, and carrying out the remaining derivation as before.
 65 This also allows to apply other adaptive methods like ADAM. On all methods, the generator and
 66 discriminator are initially chasing each other across the strategy space, producing the typical cycling
 67 pattern. When using SGA, ConOpt, or OGDA, however, eventually the algorithm diverges with the
 68 generator either mapping all the mass far away from the mode, or collapsing the generating map
 69 to become zero. Therefore, we also tried decreasing the stepsize to 0.001, which however did not
 70 prevent the divergence. For CGD, after some initial cycles the Generator starts splitting the mass
 71 and distributes it roughly evenly among the two modes. During our experiments, this configuration
 72 appeared to be robust. In the supplement, we have included a number of visualizations of the games
 73 trajectories for a variety of stepsizes and algorithms. Here, for example the folder with the name
 74 `two_mode_conOpt_25` contains the experiment with ConOpt and stepsize $\eta = 25 * 0.001 = 0.025$.

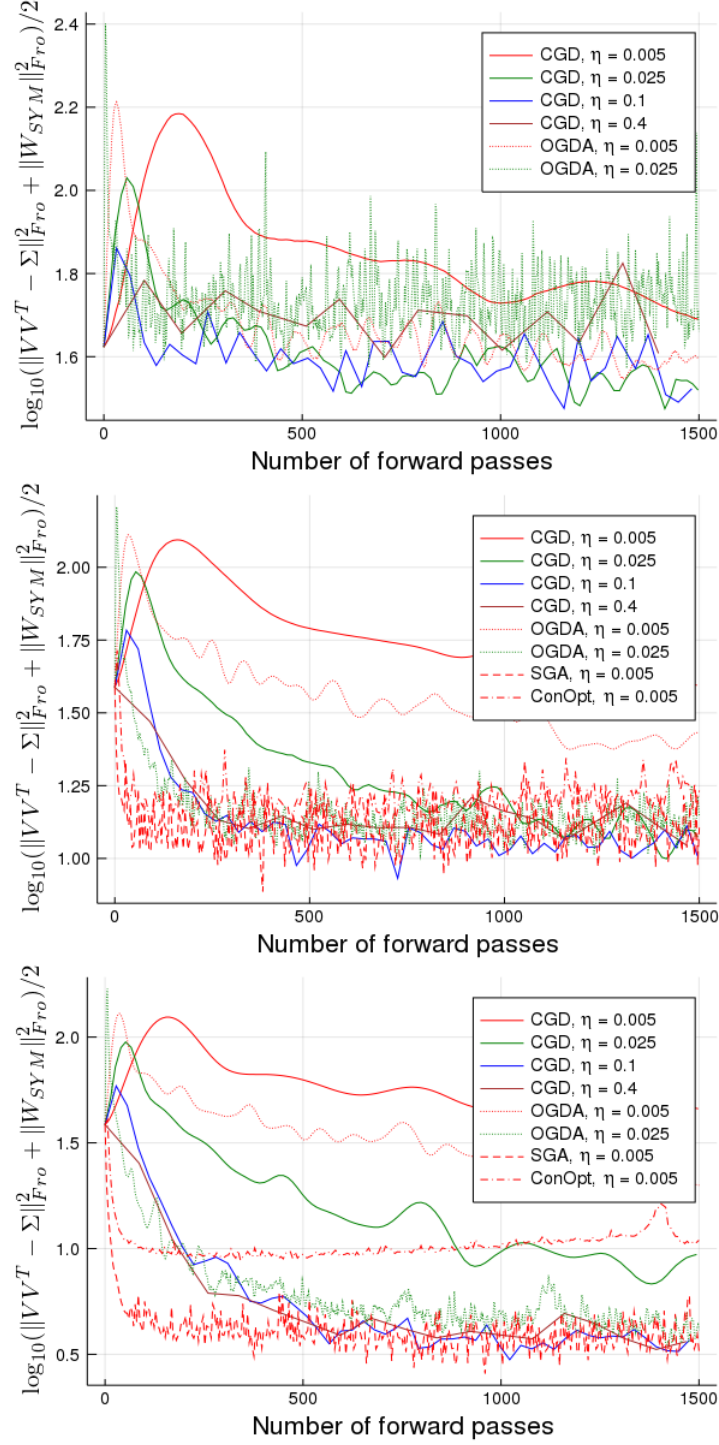


Figure 2: The decay of the residual as a function of the number of forward iterations in the stochastic case with $d = 20$ and batch sizes of 100, 1000, 10000, from top to bottom).