

GEOADALER: GEOMETRIC INSIGHTS INTO ADAPTIVE STOCHASTIC GRADIENT DESCENT ALGORITHMS

CHINEDU ELEH^{†,§}, MASUZYO MWANZA[†], EKENE AGUEGBOH[‡], HANS-WERNER VAN WYK[†]

Abstract. The Adam optimization method has achieved remarkable success in addressing contemporary challenges in stochastic optimization. This method falls within the realm of adaptive sub-gradient techniques, yet the underlying geometric principles guiding its performance have remained shrouded in mystery, and have long confounded researchers. In this paper, we introduce GeoAdaLer (Geometric Adaptive Learner), a novel adaptive learning method for stochastic gradient descent optimization, which draws from the geometric properties of the optimization landscape. Beyond emerging as a formidable contender, the proposed method extends the concept of adaptive learning by introducing a geometrically inclined approach that enhances the interpretability and effectiveness in complex optimization scenarios.

Key words. adaptive learning rate, stochastic optimization, convex optimization, large-scale optimization, machine learning

MSC codes. 90C25, 90C15, 90C30, 68Q32, 68T07

1. Introduction. Stochastic gradient descent (SGD) optimization methods [25, 29] play an important role in various scientific fields. When applied to machine learning algorithms, the objective is to adjust a set of parameters with the goal of optimizing an objective function. This optimization usually involves a series of iterative adjustments made to the parameters in each step as the algorithm progresses [14]. In the vanilla gradient descent approach, the magnitude of the gradient is the predominant annealing factor causing the algorithm to take larger steps when you are away from the optimum and smaller steps when closer to an optimum [27]. This method, however, becomes less effective near an optimal point, necessitating the selection of a smaller learning rate, which in turn affects the speed of convergence. The raw magnitude of the gradient does not always align with the optimal descent step size, thus necessitating a manually chosen learning rate. If the learning rate is too big, overshooting may occur and convergence rate is slow if the learning rate is too small. This challenge has led to the development of adaptive learning algorithms [23]. In this paper, we explore gradient descent algorithms that optimize the objective function with emphasis on the update rule. In this context, we focus on the update rule by breaking it into three components: the learning rate, the annealing factor, and the descent direction. This approach allows us to evaluate the effectiveness of an algorithm by examining the impact of its learning rate, annealing factor, and descent direction on the optimization process.

In this paper, we propose GeoAdaLer (short for Geometric Adaptive Learner), a new adaptive learning method for SGD optimization that is based on the geometric properties of the objective landscape. We use cosine of θ (where θ is the acute angle between the normal to the tangent hyperplane and the horizontal hyperplane) as an annealing factor, which takes values close to zero when the optimization is traversing points close to an optimum and close to one for points far away from an optimum.

[†]Department of Mathematics and Statistics, Auburn University, Auburn, AL 36849, USA (cae0027@auburn.edu, mzm0183@auburn.edu, hzv0008@auburn.edu).

[‡]Department of Applied Economics, Auburn University, Auburn, AL 36849, USA (esa0013@auburn.edu).

[§]Corresponding author (cae0027@auburn.edu).

43 Our method has surprising similarities to other adaptive learning methods.

44 Some of the advantages of GeoAdaLer is that it introduces a geometric approach
 45 for the annealing factor and outperforms standard SGD optimization methods due
 46 to the cosine of θ . Through both theoretical analysis and empirical observation,
 47 we identified similarities between our proposed method and other adaptive learning
 48 techniques. Our method enhances the understanding of existing algorithms and opens
 49 up more opportunities for geometric interpretability of how the algorithms traverse
 50 the objective manifold.

51 Additionally, we analyze the convergence of GeoAdaLer. We frame the optimiza-
 52 tion process as a fixed-point problem and split the analysis into deterministic and
 53 stochastic cases. Under the deterministic framework, we assume convexity and the
 54 existence of a finite optimal value. We utilize the Lipschitz continuity property of the
 55 gradient to establish upper bounds of convergence. We further employ the cocoerciv-
 56 ity and quadratic upper bound properties to establish the lower bounds of convergence
 57 [34]. Under the stochastic framework, we employ the regret function, which measures
 58 the overall difference between our method and the known optimum point, ensuring
 59 that as time tends to infinity, the regret function over time tends to zero. By ensuring
 60 that the upper bound of the regret function goes to zero, we determine the overall
 61 convergence of the stochastic method to an optimum point. Both the deterministic
 62 and stochastic analyses demonstrate the robustness of GeoAdaLer and enhance our
 63 understanding of its practical applications.

64 **2. Related Work.** In the field of adaptive stochastic gradient descent algo-
 65 rithms, we continue to see improvements, often due to the need to address shortcom-
 66 ings of previous methods. The pioneering approach is AdaGrad, which focuses on the
 67 concept of per-parameter adaptive learning rates. The foundation of AdaGrad is also
 68 the source of its limitation: the monotonic accumulation of squared gradients that
 69 could prematurely stifle learning rates [37, 8].

70 Subsequent optimizers like AdaDelta, RMSprop, and the popular Adam sought
 71 to address this issue [14, 33, 37]. AdaDelta introduced a decaying average of past
 72 squared gradients, while RMSprop utilized a similar exponential decay mechanism to
 73 limit the aggressive reduction in learning rates. Adam combined RMSprop’s adaptive
 74 learning rates with the concept of momentum for smoother updates. However, Adam
 75 has been challenged by AMSGrad which modifies the algorithm in order to imbue
 76 it with ”long term memory” which addresses issues with sub-optimal convergence
 77 under specific conditions and improves empirical performance [24]. As pointed out
 78 by [21, 31, 38], the assumptions on the hyperparameters of Adam were made before
 79 constructing the counter examples in [24]. Insights gained from these examples are
 80 invaluable for the ongoing exploration of stochastic gradient descent algorithms.

81 The literature on adaptive optimization algorithms for SGD reveals that each
 82 algorithm addresses the problems evident in its predecessors. The iterative approach,
 83 while successful in many ways, arguably emphasizes the lack of fundamental intuition
 84 regarding the dynamics of the AdaGrad family of optimizers.

85 GeoAdaLer is a novel adaptive learning method for SGD, employing the proper-
 86 ties of the objective landscape. The innovative idea behind the GeoAdaLer approach
 87 is that the acute angle between the normal to the tangent plane at x and the horizon-
 88 tal plane conveys significant curvature-related information. This information could
 89 potentially recover the power of second order methods, which are not feasible in large-
 90 scale machine learning optimization. An understanding of this geometric approach to
 91 analyze adaptive methods for SGD can shed new light on the behavior of other algo-

92 rithms in the AdaGrad family, potentially revealing the rationale for their strengths
 93 and weaknesses.

94 By discussing the geometric implications of adaptive step processes, we are able
 95 to potentially come up with optimizers that:

- 96 • **have more optimal annealing:** A geometric perspective that informs
 97 strategies to control learning rate decay more effectively, preventing prema-
 98 ture convergence or extensively slow convergence.
- 99 • **provide parameter-sensitive adaptivity:** The geometry reveals how to
 100 tailor updates for individual parameters in a more principled manner.
- 101 • **increase robustness:** Understanding the geometric implications leads to
 102 optimizers that are less dependent on sensitive hyper-parameter tuning.

103 In essence, a geometric framework promises to move beyond the reactive development
 104 pattern, allowing us to proactively design adaptive optimizers that address the core
 105 issues in the AdaGrad family with greater intuition and foresight.

106 3. Mathematical Formulation.

107 **3.1. Deterministic Optimization.** Consider the minimization of the convex
 108 objective functional $f : \mathbb{R}^n \rightarrow \mathbb{R}$ using the gradient descent algorithm. The vanilla
 109 gradient descent (GD) algorithm that maximizes the benefit of *gradient annealing* for
 110 smooth functions is

$$111 \quad (3.1) \quad x_{t+1} = x_t + \delta x_t$$

112 where $\delta x_t = -\gamma g_t$ and g_t is the gradient of the objective function at time t . This is
 113 the update step and is largely responsible for how far a step is taken in the descent
 114 direction. To see its real contribution, we decompose it further into

$$115 \quad (3.2) \quad \delta x_t = -\gamma \|g_t\| \bar{g}_t$$

116 where $\gamma > 0$ is the learning rate which is responsible for manually scaling the step
 117 size, $\|g_t\|$ is the annealing factor, and $-\bar{g}_t$ is the unit vector in the descent direction.

118 It is established that the annealing factor tends to be sub-optimal and can cause
 119 overshooting which we need to compensate for by applying a smaller learning rate,
 120 increasing convergence time [23]. To combat this issue, we propose an annealing
 121 factor based on the cosine of θ_t , where θ_t is the acute angle between the normal to the
 122 tangent hyperplane and the horizontal hyperplane at iteration t . As shown in Figure
 123 1, θ_t holds a vital information about the location of the current gradient step on the
 124 objective function. We harness this information using the cosine of θ_t and prove the
 125 following theorem.

126 **THEOREM 3.1 (Geohess).** *Let θ be the acute angle between the normal to an*
 127 *objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ which is differentiable at x . Let $\|\cdot\|$ be the norm*
 128 *induced by the inner product on \mathbb{R}^n . Then*

$$129 \quad (3.3) \quad \cos \theta = \frac{\|\nabla f(x)\|}{\sqrt{\|\nabla f(x)\|^2 + 1}}.$$

130 We call Theorem 3.1 *Geohess* since this formulation mimics the curvature information
 131 traditionally found in the full hessian matrix.

132

133 *Proof.* Let x_i be an arbitrary point in \mathbb{R}^n . Then $[-\nabla f(x_i), 1]^T$ is orthogonal
 134 to $[x - x_i, y - f(x_i)]$ where (x, y) lies on the tangent hyperplane to f at x_i . i.e.,

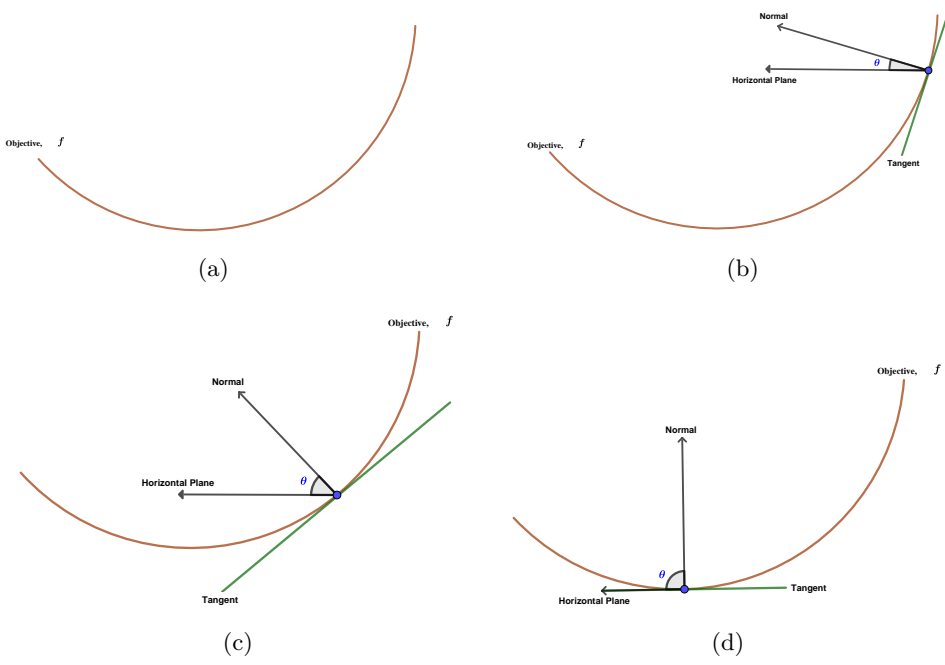


Fig. 1: Geometric view of θ as the GeoAdaLer step traverses the objective function.

135 $[-\nabla f(x_i), 1]^T$ is normal to $[x - x_i, y - f(x_i)]$. Let θ be the angle this normal makes with
 136 the horizontal hyperplane at the point $(x_i, f(x_i))$ in the descent direction $-\nabla f(x_i)$,
 137 i.e., a vector parallel to $[-\nabla f(x_i), 0]$. Then by vector calculus,

$$138 \quad (3.4) \quad \cos \theta = \frac{[-\nabla f(x_i), 1]^T \cdot [-\nabla f(x_i), 0]^T}{\|[-\nabla f(x_i), 1]^T\| \|[-\nabla f(x_i), 0]^T\|} = \frac{\|\nabla f(x_i)\|}{\sqrt{\|\nabla f(x_i)\|^2 + 1}}$$

139 as required. □

140 The proposed GeoAdaLer update step is as follows:

$$141 \quad (3.5) \quad \delta x_t = -\gamma \bar{g}_t \cos \theta_t = -\gamma \frac{\|g_t\|}{\sqrt{\|g_t\|^2 + 1}} \bar{g}_t \text{Provide}$$

142 Equation (3.5) follows from the Geohess theorem (Theorem 3.1). In equation (3.5), if
 143 g_t is zero, then the update step is zero and algorithm has converged. If g_t is not zero,
 144 then equation (3.5) reduces to

$$145 \quad \delta x_t = -\gamma \frac{g_t}{\sqrt{\|g_t\|^2 + 1}}$$

146 **3.2. Properties of GeoAdaLer Annealing Factor.** The acute angle θ_t de-
 147 pends on g_t , and therefore, the annealing factor possesses the following properties as
 148 $\|g_t\| \rightarrow 0$ and as $\|g_t\| \rightarrow \infty$ respectively:

$$149 \quad 1. \frac{\|g_t\|}{\sqrt{\|g_t\|^2 + 1}} \rightarrow 1 \quad 150 \quad 2. \frac{\|g_t\|}{\sqrt{\|g_t\|^2 + 1}} \rightarrow 0$$

151 In practice, we do not want $\|g_t\| \rightarrow \infty$ but as long as the magnitude of the gradient
 152 is large, the sufficiency in item 1 is guaranteed. Of utmost importance is item 2.

153 The GeoAdaLer algorithm is presented in Algorithm 3.1. In the deterministic
 154 setting, we set $\beta = 0$.

155 **3.3. Stochastic Optimization.** Online learning and stochastic optimization
 156 are closely linked and can be essentially used interchangeably [8, 5]. In online learning,
 157 the learner iteratively predicts a point $x_t \in X \subseteq \mathbb{R}^n$, often representing a weight vector
 158 that assigns importance values to different features. The objective of the learner is to

$$159 \quad \underset{x}{\text{minimize}} \mathbb{E}_{\xi} [f(x, \xi)],$$

160 which is too difficult to predict [10]. Instead, we minimize regret, defined as the
 161 difference between the learner’s cumulative loss and the cumulative loss of the best
 162 fixed predictor in hindsight, $x^* \in X$, chosen from a closed convex set (details Section
 163 5.2). This is performed across a sequence of convex loss functions $\{f_1, f_2, \dots\}$, where
 164 each function represents feedback or loss information available at step t [39, 30, 10].

165 Geometrically, gradients are vectors and have directions. Under a suitable distri-
 166 bution, a convex hull of these gradients approximates the true gradient for stochastic
 167 optimizations. A common and intuitive approach to approximating the expected gra-
 168 dient is by employing an exponential moving average (EMA). As is the practice in
 169 literature [12, 14], we use a momentum-like term to replace the instantaneous gradi-
 170 ent, thereby mitigating stochastic fluctuations and revealing underlying trends in the
 171 gradient values. That is,

$$172 \quad (3.6) \quad m_{t+1} = \beta m_t + (1 - \beta)g_t.$$

173 where $\beta \in [0, 1)$. We have used $g_t = \nabla f_t(x_t)$ for notational convenience. The term
 174 β is constructed as a weighted average of the historical gradients and the current
 175 gradient. This modification enhances our ability to approximate the true underlying
 176 gradient more effectively in time. It also has the ability to update and learn as new
 177 streams of data are observed.

178 In Algorithm 3.1, we present a pseudo-code of the proposed GeoAdaLer learning
 179 method.

180 **3.3.1. GeoAdaMax.** Due to stochasticity and varying frequencies of occurrence
 181 of certain model inputs, such as in deep neural networks, adaptive stochastic gradient
 182 descent methods sometimes encounter issues with non-increasing squared gradients.
 183 For instance, consider the convex objective function presented in [24] over the domain
 184 $[-1, 1]$, defined as follows:

$$185 \quad f_t(x) = \begin{cases} Cx, & \text{if } t \bmod 3 = 1 \\ -x, & \text{otherwise,} \end{cases}$$

186 where $C > 2$. In such scenarios, the adaptive step employed by the optimizer can
 187 still function effectively as a form of annealing, but the vanilla adaptive step leads to
 188 a suboptimality.

189 To address this problem, an approach was developed by [24], which involves re-
 190 taining the maximum of the normalizing denominator over iterations. This approach

Algorithm 3.1 GeoAdaLer

Require: γ : Learning rate
Require: β : Exponential decay rate for weighted average
Require: f_t : Stochastic objective function
Require: x_0 : Initial parameter vector
 $t \leftarrow 0$
while x_t not converged **do**
 $g_t \leftarrow \nabla f_t(x_t)$ (Get gradients w.r.t stochastic objective)
if $t = 0$ **then**
 $m_t \leftarrow g_t$ (initial weighted average)
else
 $m_t = \beta m_{t-1} + (1 - \beta)g_t$ (Updated weighted average)
end if
 $x_{t+1} = x_t - \gamma \cdot m_t / (\sqrt{\|m_t\|^2 + 1})$ (Update parameters)
 $t \leftarrow t + 1$
end while
return x_t

191 reduces to a self scaling SGD with momentum [24]. By implementing a similar idea, we
 192 observe related results for GeoAdaLer. We call this method GeoAdaMax, indicating
 193 the use of maximum of the variance term.

194 GeoAdaMax dynamically adjusts the denominator using the largest historical
 195 value of the EMA. When this denominator remains unchanged, the update scale re-
 196 flects its historical maximum, preserving proportionality in all future updates. This
 197 mechanism fine-tunes step sizes while maintaining alignment with past gradient mag-
 198 nitudes. The summary of the algorithm is presented in Algorithm 3.2

199 Geometrically, this is akin to increasing the angle between the normal and the
 200 horizontal plane by using a vector different from the normal. This leads to relatively
 201 smaller step sizes than if the original angle were used. Theorem 3.2 shows that indeed,
 202 the acute angle θ is increased.

203 **THEOREM 3.2.** *Let θ_t be the acute angle between the normal and the horizontal*
 204 *hyperplanes at the current iteration and let $\hat{\theta}_t$ be the acute angle between the horizontal*
 205 *hyperplane and the normal that maximizes the norm up to the current iteration, t .*
 206 *Then*

$$207 \quad \theta_t \leq \hat{\theta}_t.$$

208 *Proof.* By Algorithm 3.1 and Theorem 3.1,

$$209 \quad \cos \theta_t = \frac{\|m_t\|}{\sqrt{\|m_t\|^2 + 1}} \geq \frac{\|m_t\|}{\sqrt{\max_t \|m_t\|^2 + 1}} = \cos \hat{\theta}_t$$

210 Applying the inverse cosine function on both sides gives the desired result since \cos^{-1}
 211 is monotone decreasing on $[0, 1]$. \square

212 **4. Relationship to the AdaGrad Family.** The AdaGrad family of methods
 213 adjusts learning rates based on the accumulation of past gradients. These methods
 214 dynamically adapt the step size during optimization to handle varying gradient mag-
 215 nitudes across dimensions. For a given time step t , the learning rate adjustment in

Algorithm 3.2 GeoAdaMax

Require: γ : Learning rate
Require: β : Exponential decay rate for weighted average
Require: f_t : Stochastic objective function
Require: x_0 : Initial parameter vector
 $t \leftarrow 0$
 $u_t \leftarrow 0$
while x_t not converged **do**
 $g_t \leftarrow \nabla f_t(x_t)$ (Get gradients w.r.t stochastic objective)
if $t = 0$ **then**
 $m_t \leftarrow g_t$ (initial weighted average)
else
 $m_t = \beta m_{t-1} + (1 - \beta)g_t$ (Update weighted average)
end if
 $u_t = \max(\|m_t\|^2 + 1, u_{t-1})$
 $x_{t+1} = x_t - \gamma \cdot m_t / (\sqrt{u_t})$ (Update parameters)
 $t \leftarrow t + 1$
end while
return x_t

216 these methods depends on the accumulated gradient information, which we represent
217 as G_t . Below is a summary of how G_t is defined for the main methods in the AdaGrad
218 family:

219 **AdaGrad:** For AdaGrad, G_t is given as

$$220 \quad G_t = \sum_{i=1}^t g_i^2$$

221 where g_i is the gradient at step i . AdaGrad accumulates the squared gradients over
222 time, leading to decreasing learning rates [8].

223 **RMSProp:** The G_t step in RMSProp involves exponential moving average and
224 is given as

$$225 \quad G_t = \beta G_{t-1} + (1 - \beta)g_t^2.$$

226 RMSProp introduces an exponential decay factor β , preventing the learning rate from
227 decreasing too quickly by giving more weight to recent gradients [12].

228 AdaGrad and RMSProp modify the learning rate as follows:

$$229 \quad x_{t+1} = x_t - \frac{\gamma}{\sqrt{G_t + \epsilon}} g_t,$$

230 where η is the global learning rate and ϵ is a small constant for numerical stability.

231 **Adam:** In the case of Adam, an adjustment is not only made to the squared
232 gradients but also to the gradients where a bias-corrected moving average is applied
233 to each:

$$234 \quad m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t$$

$$235 \quad G_t = \beta_2 G_{t-1} + (1 - \beta_2)g_t^2$$

236 where m_t and G_t are the moving average of the gradients and the squared gradients
237 respective, $\hat{m}_t = m_t / (1 - \beta_1^t)$ and $\hat{G}_t = G_t / (1 - \beta_2^t)$ are the bias-corrected moving

238 average of gradients and squared gradients respectively. Adam combines the ideas of
 239 momentum and adaptive learning rates for smoother updates [14].

240 The Adam’s update rule therefore is as follows:

$$241 \quad x_{t+1} = x_t - \frac{\gamma}{\sqrt{\hat{G}_t + \epsilon}} \hat{m}_t,$$

242 **4.1. GeoAdaLer in Comparison.** The core of GeoAdaLer’s update is based
 243 on the cosine of the angle θ between the normal to the gradient and the horizontal
 244 hyperplane. It effectively uses the geometry of the optimization landscape to inform
 245 its adaptivity. Hence, with an EMA update

$$246 \quad m_t = \beta m_{t-1} + (1 - \beta)g_t,$$

247 GeoAdaLer updates as:

$$248 \quad x_{t+1} = x_t - \frac{\gamma}{\sqrt{\|m_t\|^2 + 1}} m_t.$$

249 As shown above, GeoAdaLer’s update is similar to that of the AdaGrad family with
 250 some function of the squared gradient playing a big role in the optimization step. The
 251 main differences include:

- 252 (a) The use of norm-based scaler
- 253 (b) The stability term in GeoAdaLer is naturally derived from the choice of the
 254 reference plane, for example with the choice of the normal to the gradient,
 255 our scaler produces a stability term of 1.
- 256 (c) There is only one moving average which is used in the estimation of the
 257 gradient and the function of the squared gradients
- (d) Directly comparing to Adam, we also agree that G_t is always bigger for Adam
 since by Jensen’s inequality,

$$[\beta m_{t-1} + (1 - \beta)g_t]^2 \leq \beta m_{t-1}^2 + (1 - \beta)g_t^2,$$

258 where the LHS G_t is for GeoAdaLer and the RHS for Adam.

259 We remark that all of these differences stem from the geometric intuition behind
 260 GeoAdaLer, which we believe lends itself to better understanding of the optimization
 261 paths as well as interpreting optimal values. Also, as noted in [35], the use of norm-
 262 based adaptivity ensures GeoAdaLer is robust to its hyperparameters.

263 5. Convergence Analysis.

264 **5.1. Deterministic Setting.** We analyze the convergence of GeoAdaLer, first
 265 for the deterministic case, and then for the stochastic setting. Our setup remains the
 266 same, namely:

$$267 \quad (5.1) \quad \underset{x}{\text{minimize}} \quad f(x)$$

268 using the new adaptive gradient descent (GeoAdaLer) method where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is
 269 an objective function. We recast the optimization problem (5.1) into a fixed point
 270 iteration. Let $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a nonlinear operator defined as:

$$271 \quad (5.2) \quad Tx = \left(I - \gamma \frac{\nabla f}{\sqrt{\|\nabla f\|^2 + 1}} \right) (x)$$

272 where I is the identity map.

273 THEOREM 5.1. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex and ∇f Lipschitz continuous with
 274 $\gamma \leq \frac{1}{L}$ where L is the Lipschitz constant for ∇f . Assume f attains an optimal value
 275 at $x^* = \arg \min_x f(x)$. Then T defined in (5.2) is a contraction map with contrac-
 276 tion parameter $\alpha = \sqrt{1 + \gamma^2 L_G^2 - 2\gamma \frac{L_G^2}{L}} < 1$ where L_G is the Lipschitz constant for
 277 $\frac{\nabla f}{\sqrt{\|\nabla f\|^2 + 1}}$ and $\|\cdot\|$ is the Euclidean norm in \mathbb{R}^n . That is

$$278 \quad (5.3) \quad \|Tx - Ty\| \leq \alpha \|x - y\|.$$

279 A critical component of Theorem 5.1 involves demonstrating that the mapping T
 280 is a contraction. This allows us to invoke the Banach Fixed Point Theorem, which
 281 asserts that any contraction mapping on a complete metric space possesses a unique
 282 fixed point [6, 28, 32, 11]. By iteratively applying the contraction map, starting from
 283 an initial point, we ensure convergence to this fixed point. The recursive iterations
 284 converge to the unique fixed point with a geometric rate of convergence α . This fixed
 285 point corresponds to the minimizer we are seeking [3]. We remark that Theorem 5.1
 286 holds if we replace gradient with subgradient.

287 **5.2. Stochastic Setting.** We examine the convergence properties of the GeoAda-
 288 aLer optimizer within the framework of online learning, as originally introduced in
 289 [39]. This framework involves a sequence of convex cost functions $\{f_1, f_2, \dots, f_T\}$,
 290 each of which becomes known only at its respective timestep. The objective at each
 291 step t is to estimate the parameter x_t and evaluate it using the newly revealed cost
 292 function f_t . Given the unpredictable nature of the sequence, we assess the perfor-
 293 mance of our algorithm by computing the regret. Regret is defined as the cumulative
 294 sum of the differences between the online predictions $f_t(x_t)$ and the optimal fixed
 295 parameter $f_t(x^*)$ within a feasible set X across all previous time-steps. Specifically,
 296 the regret is defined as follows [39, 14]:

$$297 \quad (5.4) \quad R(T) = \sum_{t=1}^T (f_t(x_t) - f_t(x^*))$$

298 where the optimal parameter x^* is determined by $x^* = \arg \min_{x \in X} \sum_{t=1}^T f_t(x)$. We
 299 demonstrate that GeoAdaLer achieves a regret bound of $O(\sqrt{T})$, with a detailed
 300 proof provided in the appendix C. This result aligns with the best known bounds
 301 for the general convex online learning problem. We carry over all notations from the
 302 deterministic setting. Assuming the learning rate γ_t is of order $O(t^{-1/2})$ and β_t
 303 is exponentially decaying with exponential constant λ very close to 1, we obtain the
 304 following regret bounds for online learning with GeoAdaLer algorithm.

305 THEOREM 5.2. For all $x \in \mathbb{R}^n$, and $t \leq T$, assume f_t is convex and the gradient
 306 norm $\|\nabla f_t(x)\| \leq G$. Let $\gamma_t = \frac{\gamma}{\sqrt{t}}$, $\beta_t = \beta \lambda^t$, $\lambda \in (0, 1)$, and $\beta \in [0, 1)$. For any
 307 $k \in \{1, \dots, T\}$, the separation between any point x_k generated by GeoAdaLer and
 308 the minimizer of an offline objective computed after all data is known is bounded as
 309 $\|x_k - x^*\| \leq G$. Then, for any $T \geq 1$, GeoAdaLer Algorithm achieves the regret bound:

$$310 \quad (5.5) \quad R(T) \leq \frac{D^2 \sqrt{G^2 + 1} \sqrt{T} + G(2\sqrt{T} - 1)}{2(1 - \beta)} + \frac{DG\beta(1 - \lambda^T)}{(1 - \beta)(1 - \lambda)}$$

311 From Theorem 5.2, we observe that the GeoAdaLer algorithm achieves a sublinear
 312 regret bound of $O(\sqrt{T})$ over T iterations, which is consistent with the results typically

313 found in the literature for similar algorithms. Our proof, akin to the approach taken in
 314 the Adam algorithm [14], relies significantly on the decay of β_t to ensure convergence.

315 In contrast to other adaptive stochastic gradient descent methods that adapts
 316 the current employs norm based scaling. This approach ensures that each parameter
 317 benefits from the collective dynamics of all parameters at the current time while
 318 retaining historical information in the exponential moving average of the gradients
 319 used for the adaptation. We remark that the convergence analysis for GeoAdaMax
 320 follows trivially from Theorem 5.2.

321 **COROLLARY 5.3.** *Assume that for any $x \in \mathbb{R}^n$, the function f_t is convex and*
 322 *satisfies the gradient bounds $\|\nabla f_t(x)\|_2 \leq G$. Also, assume that the distance between*
 323 *any parameter x_k generated by GeoAdaLer Algorithm and the minimizer of an offline*
 324 *objective computed after all data is known is bounded, namely $\|x_k - x^*\|_2 \leq D$ for any*
 325 *$k \in \{1, \dots, T\}$. Then, GeoAdaLer achieves the following guarantee for all $T \geq 1$:*

$$326 \quad \limsup_{T \rightarrow \infty} \frac{R(T)}{T} \leq 0$$

327 The corollary is derived by dividing the result in Theorem 5.2 by T and applying the
 328 limit superior operation to both sides of the inequality 5.5. It is important to note
 329 that when $R(T)$ yields a negative value, it indicates a favorable performance of the
 330 iterates produced by the GeoAdaLer algorithm. Specifically, such results suggest that
 331 the algorithm’s execution leads to an expected loss that is lower than that of the best
 332 possible offline algorithm, which has full foresight of all cost functions and selects for
 333 a single optimal vector as proposed in [39].

334 **6. Experiments.** We compare GeoAdaLer to other algorithms such as Adam,
 335 AMSGrad, and SGD with momentum. This comparison aims to assess how the geo-
 336 metric approach performs against these popular methods on the standard CIFAR-10
 337 and MNIST datasets. Additionally, we examine how key hyperparameters influence
 338 the convergence of GeoAdaLer across different datasets.

339 The hyperparameter settings for the algorithms are detailed as follows. These
 340 settings represent default values unless otherwise specified as recommended in the
 341 literature or considered standard for their respective packages. For SGD, the learning
 342 rate was set to 0.01, momentum to 0.9 and dampening to 0.9. For Adam & AMSGrad,
 343 the learning rate was set to 0.001, β_1 to 0.9 and β_2 to 0.99. All CPU calculations were
 344 performed on an AMD Ryzen 8-core CPU, and GPU calculations were conducted on
 345 a NVIDIA 3080ti.

346 **6.1. MNIST Dataset.** In the MNIST [18] experiment we trained a fully con-
 347 nected feed forward neural networks . It consists of three fully connected layers: the
 348 first layer takes in 784 input features (flattened 28x28 grayscale images) and outputs
 349 128 features, the second layer reduces these to 64 features, and the final layer out-
 350 puts 10 logits corresponding to class scores. Each of the first two layers is followed
 351 by a ReLU activation function. The final layer provides raw class scores. Using
 352 cross-entropy loss, GeoAdaLer and GeoAdaMax algorithms alongside the baseline
 353 optimizers were executed for 150 epochs and for 30 different weights initializations.
 354 Figure 2 and Table 1 illustrate the averaged results on the test dataset. GeoAd-
 355 aLer demonstrates comparable initial performance to algorithms such as Adam and
 356 SGD, yet it achieves better long-run performance, converging to a higher validation
 357 accuracy. GeoAdaMax further improved this performance by faster convergence.

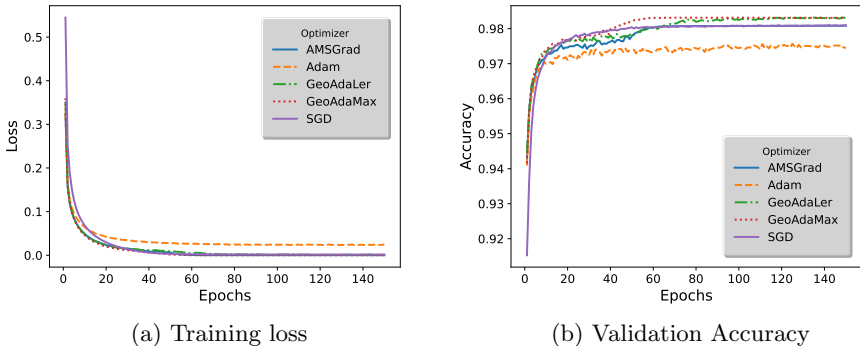


Fig. 2: MNIST

358 **6.2. CIFAR-10 Dataset.** In our CIFAR-10 [16] experiments, we run GeoAd-
 359 aLer, GoeAdaMax and the baseline optimizers for 50 epochs on a network consisting
 360 of six convolutional layers with 3×3 kernels and padding, progressively increasing in
 361 sizes of 32, 32, 64, 64, 128 and 128 filters. Each convolution is followed by a batch
 362 normalization layer and a ReLU activation. After every two convolutional layers,
 363 max-pooling with a 2×2 kernel is applied to reduce spatial dimensions, followed by
 364 dropout layers with rates 0.2, 0.3 and 0.4. The output from the convolutional block,
 365 consisting of 128 filters with a 4×4 spatial size, is flattened and passed to two fully
 366 connected layers: the first reduces the feature size to 128 with a ReLU activation
 367 and a 0.5 dropout, and the second outputs 10 logits corresponding to class scores
 368 which are passed to a softmax function. The model is trained using cross-entropy loss
 369 and re-run 30 times for each optimizer with different initializations of the weights.
 370 The averaged results on test data are shown in Figure 3 and Table 2. GeoAdaLer
 371 shows early run performance comparable to baseline optimizers and occasionally ex-
 372 ceeds them in validation accuracy. Its long run performance was only matched by
 373 Adam and GeoAdaMax. The consistent performance of GeoAdaLer in various runs
 374 reinforces the value of incorporating a geometric perspective into its design.

375 In the CIFAR-10 [16] experiment, we trained a convolutional neural network
 376 (CNN) model designed specifically for image data with RGB channels. The archi-
 377 tecture consists of six convolutional layers, progressively increasing in feature map
 378 depth (32, 64, and 128 channels), followed by max-pooling and dropout layers to re-
 379 duce overfitting and improve generalization. Each convolutional layer is followed by
 380 batch normalization to stabilize learning and expedite convergence [13]. Finally, two
 381 fully connected (linear) layers map the feature space to the 10 CIFAR-10 class scores,
 382 following common practices for classification in CNN architectures [17].

383 We train the model using cross-entropy loss for multi-class classification [9], and
 384 all adaptive gradient descent algorithms (including GeoAdaLer and GeoAdaMax)
 385 were executed for 100 epochs. To account for random initializations, we average the
 386 losses and accuracies over 30 different initializations. Model results are illustrated
 387 in Figure 3 and Table 2. GeoAdaLer demonstrates initial performance on par with
 388 other adaptive optimizers, such as Adam [14] and SGD [19], but achieved superior
 389 long-term accuracy, converging to higher validation accuracy. GeoAdaMax provided

390 further benefits, resulting in faster convergence and smoother training dynamics due
 391 to its stability near the optimal point.

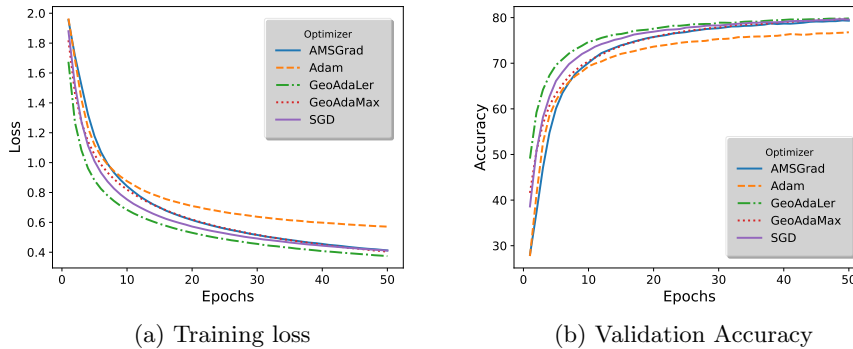


Fig. 3: CIFAR 10 Dataset

Table 1: MNIST Final Accuracy

Optimizer	Accuracy
GeoAdaLer	0.9831
GeoAdaMax	0.9831
Adam	0.9746
AMSGrad	0.9809
SGD	0.9810

Table 3: Fashion MNIST Final Accuracy

Optimizer	Accuracy
GeoAdaLer	0.9044
GeoAdaMax	0.9042
Adam	0.8838
AMSGrad	0.8993
SGD	0.8969

Table 2: CIFAR Final Accuracy

Optimizer	Accuracy
GeoAdaLer	0.7982
GeoAdaMax	0.7962
Adam	0.7679
AMSGrad	0.7932
SGD	0.7957

392 **6.3. Fashion MNIST.** In the Fashion MNIST [36] experiment we train a fully
 393 connected feed forward neural networks. It consists of three fully connected layers:
 394 the first layer takes in 784 input features (flattened 28×28 grayscale images) and
 395 outputs 512 features, the second layer reduces these to 256 features, and the final layer
 396 outputs 10 logits corresponding to class scores. Each of the first two layers is followed
 397 by a ReLU activation function. The final layer provides raw class scores. Using
 398 cross-entropy loss, GeoAdaLer and GeoAdaMax algorithms alongside the baseline

399 optimizers were executed for 100 epochs and for 30 different weights initializations.
 400 Figure 4 and Table 3 illustrate the averaged results on the test dataset. GeoAdaLer
 401 once again shows comparable results to that demonstrated by the other benchmark
 402 algorithms with GeoAdaMax showing further improvement and faster convergence on
 403 the Fashion MNIST dataset.

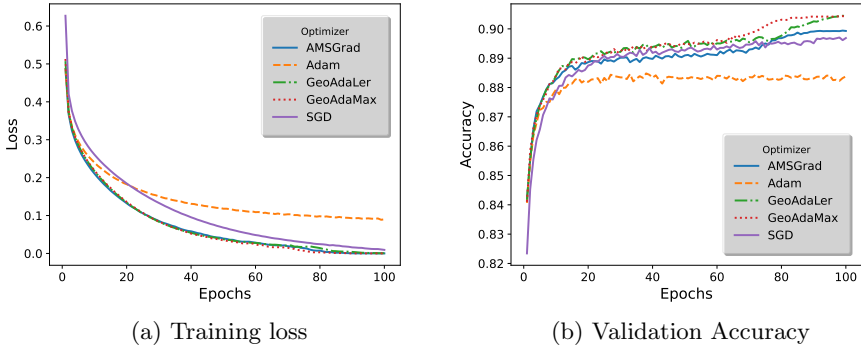


Fig. 4: Fashion MNIST Dataset

404 **6.4. Alternative Normal Plane Vectors.** By introducing a hyper-parameter
 405 ϵ in the update rule, we can explore different vectors within the normal plane:

406
$$x_{t+1} = x_t - \frac{\gamma}{\sqrt{\|m_t\|^2 + \epsilon}} m_t.$$

407 This approach allows us to select vectors with varying angles relative to the horizontal
 408 plane. This investigation stems from observing how GeoAdaMax modifies the effective
 409 angles by maximizing the denominator of the update.

410 In this experiment, we compared GeoAdaLer and GeoAdaMax on the MNIST (5)
 411 and CIFAR-10 (6) datasets using 30 different weight initializations for selected values
 412 of ϵ , all other parameters were as mentioned above in their individual experiments.
 413 The mean accuracies for each ϵ value are plotted versus the value of ϵ . The results
 414 indicate that while the normal vector may not always be the most optimal choice, the
 415 optimal ϵ value tends to be close to the default associated with the normal vector.
 416 although some improvement do exist through the use of larger values of epsilon it is
 417 dependent on the data and not constant through all the experiments.

418 **7. Conclusion.** In this paper, we investigate the adaptive stochastic gradient
 419 descent algorithm and propose a geometric approach where the normal vector to the
 420 tangent hyperplane plays a crucial role in providing curvature-like information. We
 421 call this approach GeoAdaLer, and we show that it is derived from a fundamen-
 422 tal understanding of optimization geometry. We present theoretical proof for both
 423 deterministic and stochastic settings. Empirically, we find that GeoAdaLer is com-
 424 petitively comparable with other optimization techniques. Under certain conditions,
 425 it offers better performance and stability. GeoAdaLer provides a general geometric
 426 framework applicable to most, if not all, large-scale adaptive gradient-based optimiza-
 427 tion methods. We believe that this presents a significant step towards the development
 428 of interpretable machine learning algorithms through the lens of optimization.

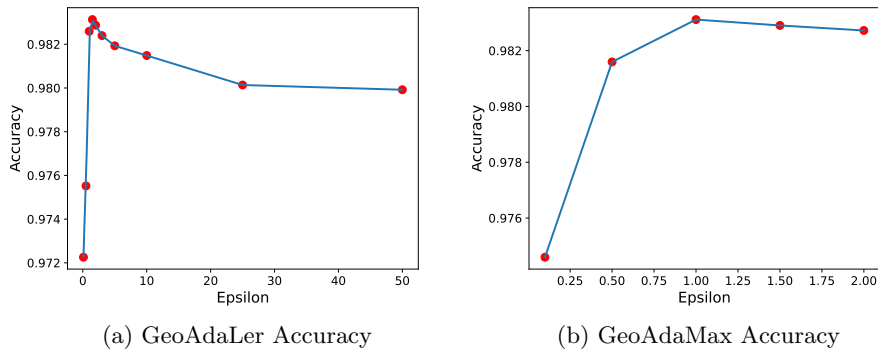


Fig. 5: MNIST Dataset

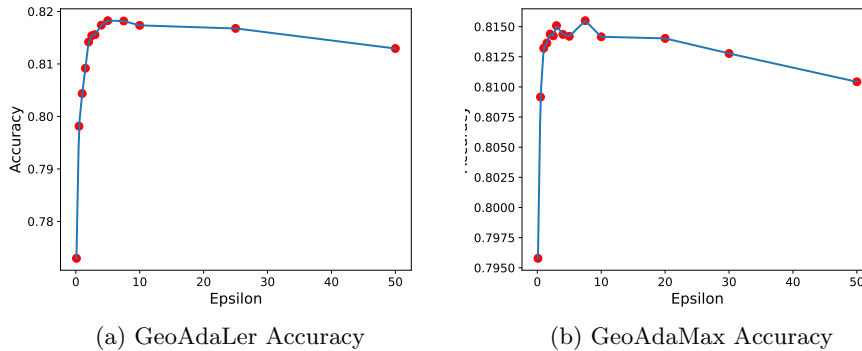


Fig. 6: CIFAR 10 Dataset

429 **8. Data and Code Availability.** Our codebase and the associated datasets
 430 used in our experiments are available in an open-source repository on GitHub:
 431 <https://github.com/Masuzyo/Geoadaler>.

432

REFERENCES

- 433 [1] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in*
 434 *Hilbert Spaces*, Springer, Cham, 2 ed., 2017, <https://doi.org/10.1007/978-3-319-48311-5>.
 435 [2] S. BOCK, J. GOPPOLD, AND M. WEISS, *An improvement of the convergence proof of the ADAM-*
 436 *optimizer*, 2018, <https://doi.org/10.48550/arXiv.1804.10587>, [https://arxiv.org/abs/1804.](https://arxiv.org/abs/1804.10587)
 437 [10587](https://arxiv.org/abs/1804.10587), <https://arxiv.org/abs/1804.10587>. arXiv preprint arXiv:1804.10587.
 438 [3] S. BOYD AND E. K. RYU, *A primer on Monotone Operator Methods (survey)*, Applied and
 439 Computational Mathematics, 15 (2016), pp. 3–43.
 440 [4] H. BRÉZIS, *Opérateurs maximaux monotones et semi-groupes de contractions dans les espaces*
 441 *de Hilbert*, vol. 5 of North-Holland Mathematics Studies, North-Holland Publishing Com-
 442 pany, Amsterdam, 1973. Notas de Matemática, Vol. 50.
 443 [5] N. CESA-BIANCHI, A. CONCONI, AND C. GENTILE, *On the generalization ability of on-line*
 444 *learning algorithms*, IEEE Transactions on Information Theory, 50 (2004), pp. 2050–2057.

- 445 [6] C. CHIDUME, *Geometric Properties of Banach Spaces and Nonlinear Iterations*, Springer,
446 Berlin, 2009.
- 447 [7] O. CHIDUME, *Foundations of mathematical real analysis: Computer science mathematical*
448 *analysis*, 2019.
- 449 [8] J. DUCHI, E. HAZAN, AND Y. SINGER, *Adaptive subgradient methods for online learning and*
450 *stochastic optimization*, Journal of Machine Learning Research, 12 (2011), pp. 2121–2159,
451 <http://jmlr.org/papers/v12/duchi11a.html>.
- 452 [9] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep Learning*, MIT Press, Cambridge, MA,
453 2016.
- 454 [10] G. J. GORDON, *Regret bounds for prediction problems*, in Proceedings of the 12th Annual
455 Conference on Computational Learning Theory (COLT), Santa Cruz, CA, USA, 1999,
456 ACM, pp. 29–40.
- 457 [11] A. GRANAS AND J. DUGUNDJI, *Fixed Point Theory*, Springer Monographs in Mathematics,
458 Springer, New York, 2003.
- 459 [12] G. HINTON, N. SRIVASTAVA, AND K. SWERSKY, *Lecture 6a: Overview of mini-batch gradient*
460 *descent*. https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf, 2012.
461 Neural Networks for Machine Learning, University of Toronto.
- 462 [13] S. IOFFE AND C. SZEGEDY, *Batch normalization: Accelerating deep network training by reducing*
463 *internal covariate shift*, in Proceedings of the 32nd International Conference on Machine
464 Learning (ICML), vol. 37 of Proceedings of Machine Learning Research, PMLR, 2015,
465 pp. 448–456, <https://proceedings.mlr.press/v37/ioffe15.html>.
- 466 [14] D. P. KINGMA AND J. BA, *Adam: A method for stochastic optimization*, 2017, <https://arxiv.org/abs/1412.6980>,
467 <https://arxiv.org/abs/1412.6980>.
- 468 [15] A. KRIZHEVSKY, *Learning multiple layers of features from tiny images*, Tech. Report
469 Technical Report, University of Toronto, 2009, [https://www.cs.toronto.edu/~kriz/](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf)
470 [learning-features-2009-TR.pdf](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf). Technical Report.
- 471 [16] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep*
472 *convolutional neural networks*, in Advances in Neural Information Processing Systems
473 Systems, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, eds., vol. 25, Curran
474 Associates, Inc., 2012, [https://proceedings.neurips.cc/paper_files/paper/2012/file/](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf)
475 [c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
- 476 [17] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convo-*
477 *lutional neural networks*, in Advances in Neural Information Processing Systems, vol. 25,
478 Curran Associates, Inc., 2012, [https://proceedings.neurips.cc/paper_files/paper/2012/file/](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf)
479 [c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
- 480 [18] Y. LECUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient-based learning applied to*
481 *document recognition*, Proceedings of the IEEE, 86 (1998), pp. 2278–2324, [https://doi.](https://doi.org/10.1109/5.726791)
482 [org/10.1109/5.726791](https://doi.org/10.1109/5.726791).
- 483 [19] Y. LECUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient-based learning applied to*
484 *document recognition*, Proceedings of the IEEE, 86 (1998), pp. 2278–2324, [https://doi.](https://doi.org/10.1109/5.726791)
485 [org/10.1109/5.726791](https://doi.org/10.1109/5.726791).
- 486 [20] Y. LECUN, C. CORTES, AND C. J. BURGES, *MNIST handwritten digit database*. [http://yann.](http://yann.lecun.com/exdb/mnist)
487 [lecun.com/exdb/mnist](http://yann.lecun.com/exdb/mnist), 2010. Accessed: 2024-05-16.
- 488 [21] H. LI, A. RAKHLIN, AND A. JADBABAIE, *Convergence of Adam under relaxed assumptions*, in
489 Advances in Neural Information Processing Systems, vol. 36, 2024. To appear.
- 490 [22] G. J. MINTY, *Monotone (nonlinear) operators in Hilbert space*, Duke Mathematical Journal,
491 29 (1962), pp. 341–346, <https://doi.org/10.1215/S0012-7094-62-02933-2>.
- 492 [23] K. NAR AND S. SASTRY, *Step size matters in deep learning*, Advances in Neural Information
493 Processing Systems, 31 (2018).
- 494 [24] S. J. REDDI, S. KALE, AND S. KUMAR, *On the convergence of adam and beyond*, in Proceedings
495 of the 6th International Conference on Learning Representations (ICLR), 2018, <https://openreview.net/forum?id=ryQu7f-RZ>.
496 [Conference version of arXiv:1904.09237](https://openreview.net/forum?id=ryQu7f-RZ).
- 497 [25] H. ROBBINS AND S. MONRO, *A Stochastic Approximation method*, The Annals of Mathematical
498 Statistics, 22 (1951), pp. 400–407, <https://doi.org/10.1214/aoms/1177729586>.
- 499 [26] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, Princeton, NJ, 1970.
- 500 [27] S. RUDER, *An overview of gradient descent optimization algorithms*, 2016, [https://doi.org/10.](https://doi.org/10.48550/arXiv.1609.04747)
501 [48550/arXiv.1609.04747](https://doi.org/10.48550/arXiv.1609.04747), <https://arxiv.org/abs/1609.04747>, [https://arxiv.org/abs/1609.](https://arxiv.org/abs/1609.04747)
502 [04747](https://arxiv.org/abs/1609.04747). arXiv preprint arXiv:1609.04747.
- 503 [28] W. RUDIN, *Principles of Mathematical Analysis*, McGraw-Hill, New York, 3rd ed., 1976.
- 504 [29] D. E. RUMELHART, G. E. HINTON, AND R. J. WILLIAMS, *Learning representations by Back-*
505 *Propagating Errors*, Nature, 323 (1986), pp. 533–536, <https://doi.org/10.1038/323533a0>.
- 506 [30] S. SHALEV-SHWARTZ, *Online learning and online convex optimization*, Foundations and Trends

- 507 in Machine Learning, 4 (2012), pp. 107–194, <https://doi.org/10.1561/22000000018>.
- 508 [31] N. SHI, D. LI, M. HONG, AND R. SUN, *RMSprop converges with proper hyper-parameter*, in
509 Proceedings of the International Conference on Learning Representations (ICLR), 2021,
510 <https://openreview.net/forum?id=3UDSdyIcBDA>. OpenReview submission.
- 511 [32] W. A. SUTHERLAND, *Introduction to Metric and Topological Spaces*, Oxford University Press,
512 Oxford, 2nd ed., 2009.
- 513 [33] T. TIELEMAN, *Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent*
514 *magnitude*, 2012, <https://cir.nii.ac.jp/crid/1370017282431050757>.
- 515 [34] L. VANDENBERGHE, *Optimization methods for large-scale systems lecture notes*, 2022, <https://www.seas.ucla.edu/~vandenbe/236C/lectures/gradient.pdf>. Lecture notes, UCLA.
- 517 [35] R. WARD, X. WU, AND L. BOTTOU, *AdaGrad stepsizes: Sharp convergence over noncon-*
518 *convex landscapes*, in Proceedings of the 36th International Conference on Machine Learning,
519 K. Chaudhuri and R. Salakhutdinov, eds., vol. 97 of Proceedings of Machine Learning
520 Research, PMLR, 09–15 Jun 2019, pp. 6677–6686, <https://proceedings.mlr.press/v97/ward19a.html>.
- 522 [36] H. XIAO, K. RASUL, AND R. VOLLGRAF, *Fashion-MNIST: A novel image dataset for bench-*
523 *marking machine learning algorithms*, 2017, <https://arxiv.org/abs/1708.07747>, <https://arxiv.org/abs/1708.07747>. arXiv preprint arXiv:1708.07747.
- 525 [37] M. D. ZEILER, *ADADELTA: An Adaptive Learning Rate Method*, arXiv e-prints, (2012),
526 arXiv:1212.5701, p. arXiv:1212.5701, <https://doi.org/10.48550/arXiv.1212.5701>, <https://arxiv.org/abs/1212.5701>.
- 528 [38] Y. ZHANG, C. CHEN, N. SHI, R. SUN, AND Z.-Q. LUO, *Adam can converge without any mod-*
529 *ification on update rules*, in Advances in Neural Information Processing Systems, vol. 35,
530 2022, pp. 28386–28399.
- 531 [39] M. ZINKEVICH, *Online convex programming and generalized infinitesimal gradient ascent*, in
532 Proceedings of the 20th International Conference on Machine Learning (ICML), ICML’03,
533 Washington, DC, USA, 2003, AAAI Press, pp. 928–935, <https://dl.acm.org/doi/10.5555/3041838.3041955>.
- 534

535 **Appendix A. Deterministic Convergence Proof.**

536 DEFINITION A.1. Let $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and let $\beta > 0$. The operator G is said to be
537 β -cocoercive if

538 (A.1) $\langle G(x) - G(y), x - y \rangle \geq \beta \|G(x) - G(y)\|^2$ for all $x, y \in \mathbb{R}^n$.

539 Equivalently, for $L := 1/\beta > 0$, G is $(1/L)$ -cocoercive if

540 (A.2) $\frac{1}{L} \|G(x) - G(y)\|^2 \leq \langle G(x) - G(y), x - y \rangle$ for all $x, y \in \mathbb{R}^n$.

541 When $G = \nabla f$ for a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we say that f has β -
542 cocoercive gradient if ∇f is β -cocoercive in the sense of (A.1).

LEMMA A.2. Suppose ∇f is Lipschitz continuous with parameter L , domain of f is \mathbb{R}^n and f has a minimum at x^* . Then

$$\frac{1}{2L} \|\nabla f(z)\|^2 \leq f(z) - f(x^*).$$

543 *Proof.* The proof relies on the following quadratic upper bound property

544 (A.3) $f(y) \leq f(z) + \langle \nabla f(z), y - z \rangle + \frac{L}{2} \|y - z\|^2$ for all $y, z \in \text{dom}(f)$

545 where $\text{dom}(f)$ is a convex set. Taking infimum on both sides of A.3 gives

546
$$f(x^*) = \inf_y f(y) \leq \inf_y \left(f(z) + \langle \nabla f(z), y - z \rangle + \frac{L}{2} \|y - z\|^2 \right)$$

547
$$= \inf_{\|v\|=1} \inf_t \left(f(z) + t \langle \nabla f(z), v \rangle + \frac{Lt^2}{2} \right)$$

548
$$= \inf_{\|v\|=1} \left(f(z) - \frac{1}{2L} \langle \nabla f(z), v \rangle \right)$$

549 (A.4)
$$= f(z) - \frac{1}{2L} \|\nabla f(z)\|^2.$$

550 Rearranging gives the claim. □

551 LEMMA A.3 (Cocoercivity). Assume f is convex, proper and lower semicontinuous.
552 Let ∇f be Lipschitz and let L_G be the Lipschitz constant for $\frac{\nabla f}{\sqrt{\|\nabla f\|^2 + 1}}$.
553 Then the following cocoercivity property holds:

554
$$\frac{1}{L_G} \left\| \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} \right\|^2$$

555
$$\leq \left\langle \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}}, x - y \right\rangle$$

556 *Proof.* For any x and y , let

557
$$f_x(z) := F(z) - \frac{\nabla f(x)^T z}{\sqrt{\|\nabla f(x)\|^2 + 1}},$$

558
$$f_y(z) := F(z) - \frac{\nabla f(y)^T z}{\sqrt{\|\nabla f(y)\|^2 + 1}}$$

559 for some constant a where $F(z) := \int_a^z \frac{\nabla f(u)}{\sqrt{\|\nabla f(u)\|^2 + 1}} du$ is a scalar potential func-
 560 tion for the vector field integrand. Notice that F is well defined since the integrand
 561 is a conservative vector field. Also, $F(z)$ is convex since the integrand is a maxi-
 562 mal monotone operator resulting from a convex, proper and lower semi-continuous
 563 function, f [22, 4, 26].

564 Thus, f_x and f_y are well defined and convex since the difference of a convex
 565 function and a linear function is convex. f_x is minimized at $z = x$, thus evaluating
 566 f_x at y and x and subtracting the results gives

$$\begin{aligned}
 567 \quad & F(y) - F(x) - \left\langle \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}}, y - x \right\rangle \\
 568 \quad & = f_x(y) - f_x(x) \\
 569 \quad (A.5) \quad & \geq \frac{1}{2L_G} \|\nabla f_x(y)\|^2 \\
 570 \quad (A.6) \quad & = \frac{1}{2L_G} \left\| \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} - \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} \right\|^2.
 \end{aligned}$$

571 Similarly, $z = y$ minimizes f_y and

$$\begin{aligned}
 572 \quad & F(x) - F(y) - \left\langle \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}}, x - y \right\rangle \\
 573 \quad & = f_y(x) - f_y(y) \\
 574 \quad & \geq \frac{1}{2L_G} \|\nabla f_y(x)\|^2 \\
 575 \quad (A.7) \quad & = \frac{1}{2L_G} \left\| \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} - \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} \right\|^2.
 \end{aligned}$$

576 Adding A.6 and A.7 yields the result. \square

LEMMA A.4 (Existence of $\text{Fix}(T)$). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be proper, lower semicon-
 tinuous and convex. If f is coercive, i.e.,*

$$\|x\| \rightarrow \infty \text{ implies } f(x) \rightarrow \infty,$$

then f attains a minimizer $x^ \in \arg \min_x f(x)$ and $\arg \min_x f \neq \emptyset$. If in addition f is
 differentiable, then*

$$\arg \min_x f = \{x \in \mathbb{R}^n : \nabla f(x) = 0\}.$$

Consequently, for the operator T defined in (5.2),

$$\text{Fix}(T) = \arg \min_x f \neq \emptyset.$$

Proof. Existence of a minimizer for proper, lower semicontinuous, coercive convex
 functions on \mathbb{R}^n is standard; see, e.g., [1, 26]. For differentiable convex f , first-order
 optimality gives $\nabla f(x) = 0$ if and only if $x \in \arg \min_x f$ (see, e.g., [1, Proposition 17.6]).

Finally, since

$$T(x) = x - \gamma \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}},$$

577 we have $T(x) = x$ if and only if $\nabla f(x) = 0$, hence $\text{Fix}(T) = \arg \min_x f$. \square

THEOREM A.5 (Deterministic convergence). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex and continuously differentiable. Assume in addition that f is proper and lower semicontinuous, and that f is coercive. Let*

$$g(x) := \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} \quad \text{and} \quad T(x) := (I - \gamma g)(x),$$

where I is the identity map and $L_G > 0$ denotes the Lipschitz constant for g . Then f attains a minimizer $x^* \in \arg \min f$, and $\text{Fix}(T) = \arg \min_x f \neq \emptyset$. If $0 < \gamma < \frac{2}{L_G}$, then T is nonexpansive,

$$\|Tx - Ty\| \leq \|x - y\| \quad \text{for all } x, y \in \mathbb{R}^n.$$

Moreover, for the fixed-point iteration $x_{k+1} = T(x_k)$, the sequence $\{x_k\}$ converges to a point $x^* \in \text{Fix}(T) = \arg \min_x f$.

Proof. First, we compute as follows

$$\begin{aligned} \|Tx - Ty\|^2 &= \left\| x - \frac{\gamma \nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - y + \frac{\gamma \nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} \right\|^2 \\ &= \left\| (x - y) - \gamma \left(\frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} \right) \right\|^2 \\ &= \|x - y\|^2 - 2\gamma \left\langle x - y, \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} \right\rangle \\ &\quad + \gamma^2 \left\| \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} \right\|^2 \\ &\leq \|x - y\|^2 + \left(\gamma^2 - \frac{2\gamma}{L_G} \right) \left\| \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} \right\|^2 \\ &= \|x - y\|^2 - \gamma \left(\frac{2}{L_G} - \gamma \right) \left\| \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} - \frac{\nabla f(y)}{\sqrt{\|\nabla f(y)\|^2 + 1}} \right\|^2. \end{aligned}$$

Inequality (A.8) follows from cocoersivity (Lemma A.3). If $0 < \gamma < 2/L_G$, then (A.9) implies

$$\|Tx - Ty\|^2 \leq \|x - y\|^2,$$

hence T is nonexpansive.

By Lemma A.4, f attains a minimizer $x^* \in \arg \min_x f$. Hence, $\text{Fix}(T) = \arg \min_x f \neq \emptyset$. Choose $x^* \in \text{Fix}(T)$ and take $y = x^*$ in (A.9). Since $T(x^*) = x^*$, we have $\nabla f(x^*) = 0$ and thus

$$\|x_{k+1} - x^*\|^2 \leq \|x_k - x^*\|^2 - \gamma \left(\frac{2}{L_G} - \gamma \right) \left\| \frac{\nabla f(x_k)}{\sqrt{\|\nabla f(x_k)\|^2 + 1}} \right\|^2. \quad \square$$

Thus $\{\|x_k - x^*\|\}$ is nonincreasing and $\{x_k\}$ is bounded.

Summing (A.10) from $k = 0$ to N yields

$$\gamma \left(\frac{2}{L_G} - \gamma \right) \sum_{k=0}^N \left\| \frac{\nabla f(x_k)}{\sqrt{\|\nabla f(x_k)\|^2 + 1}} \right\|^2 \leq \|x_0 - x^*\|^2,$$

so

$$\sum_{k=0}^{\infty} \left\| \frac{\nabla f(x_k)}{\sqrt{\|\nabla f(x_k)\|^2 + 1}} \right\|^2 < \infty.$$

In particular,

$$\left\| \frac{\nabla f(x)}{\sqrt{\|\nabla f(x)\|^2 + 1}} \right\| \text{ vanishes iff } \|\nabla f(x)\| = 0$$

593 Since $\|x_{k+1} - x_k\| = \gamma \left\| \frac{\nabla f(x_k)}{\sqrt{\|\nabla f(x_k)\|^2 + 1}} \right\|$, it follows that $\|x_{k+1} - x_k\| \rightarrow 0$ as $k \rightarrow \infty$.

594 By boundedness of $\{x_k\}$, the Bolzano-Weierstrass theorem implies existence of a con-
595 vergent subsequence, namely there exist $\{k_j\}$ and $\bar{x} \in \mathbb{R}^n$ such that $x_{k_j} \rightarrow \bar{x}$ [1, 7].
596 By continuity of ∇f and $\nabla f(x_{k_j}) \rightarrow 0$, we obtain $\nabla f(\bar{x}) = 0$, hence $\bar{x} \in \text{Fix}(T)$.

597

598 Finally, (A.10) implies that for this \bar{x} the sequence $\|x_k - \bar{x}\|^2$ is nonincreasing and
599 has a subsequence converging to 0. Hence $\|x_k - \bar{x}\|^2 \rightarrow 0$. Therefore $x_k \rightarrow \bar{x} \in \text{Fix}(T)$.

600

A best-iterate residual rate follows directly from the telescoping bound. For any $N \geq 0$,

$$\min_{0 \leq k \leq N} \left\| \frac{\nabla f(x_k)}{\sqrt{\|\nabla f(x_k)\|^2 + 1}} \right\|^2 \leq \frac{\|x_0 - x^*\|^2}{\gamma(2/L_G - \gamma)(N + 1)}.$$

601 Appendix B. Stochastic Convergence Proof.

602 **B.1. Important Lemmas and Definitions.** In this section, we provide proof
603 of convergence of our algorithm. To this end, we start with some definitions and
604 lemmas necessary for the main theorem.

605 DEFINITION B.1. A differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for all $x, y \in$
606 \mathbb{R} ,

$$607 \text{ (B.1)} \quad f(y) \geq f(x) + \nabla f(x)^T (y - x).$$

608 For the rest of the paper, we make the following assumptions on the stochastic objec-
609 tive function $f_t : \mathbb{R}^n \rightarrow \mathbb{R}$ where the iteration counter $t \geq 1$.

- 610 ASSUMPTIONS B.2. 1. f_t is convex for all t .
611 2. For all t , f_t is differentiable.
612 3. For all t , there exists $G \geq 0$ such that $\|\nabla f_t(x)\| \leq G$ for all $x \in X \subseteq \mathbb{R}^n$
613 where X is the feasible set.
614 4. $A := \{x_1, x_2, \dots\}$, the iterates generated by GeoAdaLer algorithm.
615 5. $x^* := \arg \min_{x \in X} \sum_{t=1}^T f_t(x)$ exists and $\|x_k - x^*\| \leq D$ for all $x_k \in A$.
616 6. $\beta_t := \beta \lambda^{t-1}$ where $\lambda \in (0, 1), \beta \in [0, 1)$.

617 Again, for notational convenience, we take $g_t = \nabla f_t(x_t)$.

618 LEMMA B.3. Under Assumptions B.2, no. 3, the exponential moving average

$$619 \text{ (B.2)} \quad m_t = \beta_t m_{t-1} + (1 - \beta_t) g_t$$

620 is bounded for all t .

621 *Proof.* Iteratively expanding out m_t , we obtain

$$622 \quad m_t = (1 - \beta_t) \sum_{i=1}^t \beta_t^{t-i+1} g_i.$$

623 Taking the Euclidean norm on both sides and using the boundedness of g_t gives

$$\begin{aligned} 624 \quad \|m_t\| &= \|(1 - \beta_t) \sum_{i=1}^t \beta_t^{t-i+1} g_i\| \\ 625 &\leq (1 - \beta_t) \sum_{i=1}^t \beta_t^{t-i+1} \|g_i\| \\ 626 &\leq (1 - \beta_t) G \sum_{i=1}^t \beta_t^{t-i+1} \\ 627 &= G(1 - \beta_t^t) \leq G \end{aligned}$$

628

□

629 LEMMA B.4. *Under Assumptions B.2, no. 6, the following inequalities hold*

- 630 1. $\frac{\beta_t}{1 - \beta_t} \leq \frac{\beta}{1 - \beta}$
- 631 2. $\frac{1}{1 - \beta_t} \leq \frac{1}{1 - \beta}$
- 632 3. $\sum_{t=1}^T \frac{\beta_t}{1 - \beta_t} \leq \frac{\beta(1 - \lambda^T)}{(1 - \beta)(1 - \lambda)}$

633 for all t .

634 *Proof.* 1). For all t , $\lambda < 1$ gives $\lambda^{t-1} < 1$ so that $\frac{1}{\lambda^{1-t}} < 1$. Thus

$$635 \quad \frac{\beta_t}{1 - \beta_t} = \frac{\beta \lambda^{t-1}}{1 - \beta \lambda^{t-1}} = \frac{\beta}{\lambda^{1-t} - \beta} \leq \frac{\beta}{1 - \beta}.$$

2). $\beta_t = \beta \lambda^{t-1} \leq \beta$ for all t since $\lambda \in (0, 1)$, $\beta \in [0, 1)$. So, $1 - \beta \leq 1 - \beta \lambda^{t-1}$ leads to

$$\frac{1}{1 - \beta_t} \leq \frac{1}{1 - \beta},$$

636 as required.

637 3). By Lemma B.4 no. 2, we have

$$638 \quad \sum_{t=1}^T \frac{\beta_t}{1 - \beta_t} \leq \sum_{t=1}^T \frac{\beta_t}{1 - \beta} = \frac{\beta}{1 - \beta} \sum_{t=1}^T \lambda^{t-1} = \frac{\beta(1 - \lambda^T)}{(1 - \beta)(1 - \lambda)}. \quad \square$$

639 LEMMA B.5. *For all $t \geq 1$ and for all $T \geq t$, the following inequality holds*

$$640 \quad (\text{B.3}) \quad \sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 1 + \int_1^T \frac{1}{\sqrt{t}} dt.$$

641 Proof trivially follows from the integral test for convergence of series and is also given
642 in [2].

643 LEMMA B.6. *Under Assumptions B.2, no. 1 and 2, the following inequality holds*
 644 *for all t*

$$645 \quad (B.4) \quad R(T) \leq \sum_{t=1}^T g_t \cdot (x_t - x^*).$$

646 *Proof.* By convexity of f_t for each t ,

$$647 \quad f_t(x^*) - f_t(x_t) \geq \nabla f_t(x_t) \cdot (x^* - x_t).$$

648 It then follows that

$$649 \quad f_t(x_t) - f_t(x^*) \leq \nabla f_t(x_t) \cdot (x_t - x^*) = g_t \cdot (x_t - x^*).$$

650 Hence, summing both sides from $t = 1$ to T gives the required result. \square

651 **Appendix C. Proof of Theorem 5.2.** From the update rule in Algorithm
 652 3.1

$$653 \quad x_{t+1} = x_t - \gamma_t \frac{m_t}{\sqrt{\|m_t\|^2 + 1}}.$$

654 Subtracting x^* from both sides and taking norm squared, we obtain

$$655 \quad \|x_{t+1} - x^*\|^2 = \left| x_t - \gamma_t \frac{m_t}{\sqrt{\|m_t\|^2 + 1}} \right|^2 \\ 656 \quad = \|x_t - x^*\|^2 - \frac{2\gamma_t}{\sqrt{\|m_t\|^2 + 1}} m_t \cdot (x_t - x^*) + \gamma_t^2 \frac{\|m_t\|^2}{\|m_t\|^2 + 1}.$$

657 Substituting $m_t = \beta_t m_{t-1} + (1 - \beta_t) g_t$, we obtain

$$658 \quad \|x_{t+1} - x^*\|^2 = \|x_t - x^*\|^2 - \frac{2\gamma_t}{\sqrt{\|m_t\|^2 + 1}} m_{t-1} \cdot (x_t - x^*) \\ - \frac{2\gamma_t(1 - \beta_t)}{\sqrt{\|g_t\|^2 + 1}} g_t \cdot (x_t - x^*) + \gamma_t^2 \frac{\|m_t\|^2}{\|m_t\|^2 + 1}.$$

659 Rearrange to have $g_t \cdot (x_t - x^*)$ on the left hand side:

$$660 \quad g_t \cdot (x_t - x^*) = \frac{\sqrt{\|m_t\|^2 + 1}}{2\gamma_t(1 - \beta_t)} [\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] - \frac{\beta_t}{1 - \beta_t} m_{t-1} \cdot (x_t - x^*) \\ + \frac{\gamma_t}{2(1 - \beta_t)} \frac{\|m_t\|^2}{\sqrt{\|m_t\|^2 + 1}} \\ 661 \quad (C.1) \quad \leq \frac{\sqrt{\|m_t\|^2 + 1}}{2\gamma_t(1 - \beta_t)} [\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] \\ + \frac{\beta_t}{1 - \beta_t} \|m_{t-1}\| \|x^* - x_t\| + \frac{\gamma_t \|m_t\|^2}{2(1 - \beta_t)}.$$

662 where inequality C.1 follows from the Cauchy-Schwartz inequality applied to the sec-
 663 ond term $-m_{t-1} \cdot (x_t - x^*) = m_{t-1} \cdot (x^* - x_t)$ and the fact that $\frac{\|m_t\|^2}{\sqrt{\|m_t\|^2 + 1}} \leq 1$ for all
 664 t .

665 Summing both sides from $t = 1$ to T and using Lemma B.3 and Assumptions B.2 no.
 666 3 and 5, we further obtain

$$\begin{aligned}
 667 \quad \sum_{t=1}^T g_t \cdot (x_t - x^*) &\leq \sum_{t=1}^T \frac{\sqrt{G^2 + 1}}{2\gamma_t(1 - \beta_t)} [\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] \\
 &\quad + DG \sum_{t=1}^T \frac{\beta_t}{1 - \beta_t} + G^2 \sum_{t=1}^T \frac{\gamma_t}{2(1 - \beta_t)} \\
 668 \quad (\text{C.2}) \quad &\leq \sum_{t=1}^T \frac{\sqrt{G^2 + 1}}{2\gamma_t(1 - \beta_t)} [\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] \\
 &\quad + \frac{DG\beta(1 - \lambda^T)}{(1 - \beta)(1 - \lambda)} + \frac{G^2}{2(1 - \beta)} \sum_{t=1}^T \gamma_t.
 \end{aligned}$$

669 where the last inequality follows from Assumptions B.2 no. 3 and Lemma B.4 no. 2
 670 and 3.

671 By expanding the first term and and rewriting the what is left in compact form, we
 672 obtain

$$\begin{aligned}
 673 \quad \sum_{t=1}^T \frac{1}{\gamma_t(1 - \beta_t)} [\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] \\
 674 \quad = \frac{1}{\gamma_1(1 - \beta_1)} \|x_1 - x^*\|^2 - \frac{1}{\gamma_T(1 - \beta_T)} \|x_{T+1} - x^*\|^2 \\
 \quad + \sum_{t=2}^T \left(\frac{1}{\gamma_t(1 - \beta_t)} - \frac{1}{\gamma_{t-1}(1 - \beta_{t-1})} \right) \|x_t - x^*\|^2 \\
 675 \quad \leq \frac{\|x_1 - x^*\|^2}{\gamma_1(1 - \beta_1)} \\
 \quad + \sum_{t=1}^T \left(\frac{1}{\gamma_t(1 - \beta_t)} - \frac{1}{\gamma_{t-1}(1 - \beta_{t-1})} \right) \|x_t - x^*\|^2.
 \end{aligned}$$

676

677 where the last inequality follows from dropping the second term. By the Assumptions
 678 B.2 no. 5, it further simplifies to

$$\begin{aligned}
 679 \quad \sum_{t=1}^T \frac{1}{\gamma_t(1 - \beta_t)} [\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] \\
 680 \quad \leq \frac{D^2}{\gamma_1(1 - \beta_1)} + D^2 \sum_{t=2}^T \left(\frac{1}{\gamma_t(1 - \beta_t)} - \frac{1}{\gamma_{t-1}(1 - \beta_{t-1})} \right) \\
 681 \quad = \frac{D^2}{\gamma_1(1 - \beta_1)} - \frac{D^2}{\gamma_1(1 - \beta_1)} + \frac{D^2}{\gamma_T(1 - \beta_T)} \\
 682 \quad = \frac{D^2}{\gamma_T(1 - \beta_T)}.
 \end{aligned}$$

From Lemma B.4 we have.

$$\sum_{t=1}^T \frac{1}{\gamma_t(1-\beta_t)} [\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] \leq \frac{D^2}{\gamma_T(1-\beta)}.$$

683 Therefore equation C.2 become:

$$684 \quad (C.3) \quad \sum_{t=1}^T g_t \cdot (x_t - x^*) \leq \frac{D^2\sqrt{G^2+1}}{\gamma_T(1-\beta_T)} + \frac{DG\beta(1-\lambda^T)}{(1-\beta)(1-\lambda)} + \frac{G^2}{2(1-\beta)} \sum_{t=1}^T \gamma_t.$$

685 Assuming $\gamma_t = \frac{1}{\sqrt{t}}$ we further obtain:

$$686 \quad (C.4) \quad \sum_{t=1}^T g_t \cdot (x_t - x^*) \leq \frac{\sqrt{T}D^2\sqrt{G^2+1}}{(1-\beta_T)} + \frac{DG\beta(1-\lambda^T)}{(1-\beta)(1-\lambda)} + \frac{G^2}{2(1-\beta)} \sum_{t=1}^T \frac{1}{\sqrt{t}}$$

$$687 \quad (C.5) \quad \leq \frac{\sqrt{T}D^2\sqrt{G^2+1}}{(1-\beta_T)} + \frac{DG\beta(1-\lambda^T)}{(1-\beta)(1-\lambda)} + \frac{G^2(2T-1)}{2(1-\beta)}.$$

Equation C.5 follows from Lemma B.5 and so our regret is bounded above by:

$$R(T) \leq \frac{\sqrt{T}D^2\sqrt{G^2+1}}{(1-\beta_T)} + \frac{DG\beta(1-\lambda^T)}{(1-\beta)(1-\lambda)} + \frac{G^2(2T-1)}{2(1-\beta)}.$$

688 **Appendix D. Datasets . MNIST:** The MNIST database of handwritten digits.
689 Licensed under the Creative Commons Attribution 4.0 License[20].

690 **CIFAR-10:** The CIFAR-10 dataset consists of 60000 32x32 colour images in 10
691 classes, with 6000 images per class. There are 50000 training images and 10000 test
692 images. Licensed under the Creative Commons Attribution 4.0 License [15].

693 **Fashion MNIST:** The Fashion MNIST database of fashion images. Licensed
694 under The MIT License (MIT).[36]