EXACT LINEAR-RATE GRADIENT DESCENT: OPTIMAL ADAPTIVE STEPSIZE THEORY AND PRACTICAL USE

Anonymous authors

Paper under double-blind review

Abstract

Consider gradient descent iterations $x^{k+1} = x^k - \alpha_k \nabla f(x^k)$. Suppose gradient exists and $\nabla f(x^k) \neq 0$. We propose the following closed-form stepsize choice:

$$\alpha_k^{\star} = \frac{\|\boldsymbol{x}^{\star} - \boldsymbol{x}^k\|}{\|\nabla f(\boldsymbol{x}^k)\|} \cos \eta_k, \qquad (\text{theoretical})$$

where η_k is the angle between vectors $\boldsymbol{x}^* - \boldsymbol{x}^k$ and $-\nabla f(\boldsymbol{x}^k)$. It is universally applicable and admits an exact linear convergence rate with factor $\sin^2 \eta_k$. Moreover, if f is convex and L-smooth, then $\alpha_k^* \ge 1/L$.

For practical use, we approximate (can be exact) the above via

$$\alpha_k^{\dagger} = \gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - f_0}{\|\nabla f(\boldsymbol{x}^k)\|^2}, \qquad (\text{practical use})$$

where γ_0 is a tunable parameter; \bar{f}_0 is a guess on the smallest objective value (can be auto. updated). Suppose f is convex and $\bar{f}_0 = f(x^*)$, then any choice from $\gamma_0 \in (0, 2]$ guarantees an exact linear-rate convergence to the optimal point.

We consider a few examples. (i) An \mathbb{R}^2 quadratic program, where a well-known ill-conditioning bottleneck is addressed, with a rate strictly better than $O(1/2^k)$. (ii) A geometric program, where an inaccurate guess \overline{f}_0 remains powerful. (iii) A non-convex MNIST classification problem via neural networks, where preliminary tests show that ours admits better performance than the state-of-the-art algorithms, particularly a tune-free version is available in some settings.

031 032

004

010

011

016 017

018

024

025

026

027

028

029

033 1 INTRODUCTION

The gradient descent (GD) algorithm, dated back to Cauchy in 1847, is arguably the most popular iterative algorithm. It is often treated as the default optimizer for neural networks Rumelhart et al. (1986); Ruder (2016); Goodfellow et al. (2016). GD's procedure is remarkably simple: repeatedly subtract the current iterate with its gradient. However, such a raw version suffers from a serious issue — it almost always overshoots the minimum. To guarantee convergence, damping the gradient by a stepsize α is necessary. How to properly choose such a stepsize is one of the most headache issues, since a large choice would overshoot and a small one leads to slow convergence. In practice, the stepsize (a.k.a. learning rate) is "often the single most important hyper-parameter" Bengio (2012).

043 To our best knowledge, in the current literature, a general convergence guarantee for GD only exists 044 in the convex case, and requires at least one strong assumption, the L-smoothness. Specifically, if 045 one can access the Lipschitz constant L, then any choice from $\alpha \in (0, 2/L)$ guarantees convergence, with 1/L the default choice, see e.g. (Ryu & Yin, 2022, Sec. 2.4.3). Despite such a guarantee being 046 available, it is rarely used directly in large-scale problems, due to L is either not computable or simply 047 too expensive. There does exist some work that allow estimation of L, see e.g., Anil et al. (2019); 048 Fazlyab et al. (2019); Combettes & Pesquet (2020). However, their focus is often not regarding the 049 stepsize selection issue, appears related to the complication of the estimation scheme and that the estimation error in L will propagate to the GD algorithm. In this manuscript, such an issue will be 051 avoided, since our result does not rely on L. 052

053 One critique of the above classical theory is that the stepsize is fixed throughout all iterations of GD. This eliminates the possibility of some large feasible stepsize choices in the middle steps and

consequently slows down the algorithm. A better strategy should be adaptively adjusting the stepsize 055 according to the current progress. Such an idea is old, at least traced back to Almeida et al. (1999). 056 The real issue is how to adjust the stepsize adaptively? In the literature, several outstanding heuristic 057 methods have been proposed, e.g., AdaGrad Duchi et al. (2011), RMSProp Tieleman & Hinton. 058 (2012), Adam Kingma & Ba (2015). However, an adaptive stepsize theory has not been established. This manuscript will fill in this blank space. In the convex case, we show the feasible stepsize selection range that guarantees convergence being $(0, 2\alpha_k^*)$, with α_k^* the optimal k-th choice. Moreover, α_k^* is 060 lower bounded by 1/L, implying the new range enlarges the aforementioned classical one (0, 2/L). 061 Also, our optimal stepsize yields an exact linear rate with factor $\sin^2 \eta_k$. Let us note that if $\sin \eta_k = 0$, 062 then GD will converge instantly, see an example in Section 4.1.2. 063

Remarkably, our theory also applies to a non-convex function. A notable difference is that the optimal choice α_k^* can be negative now, and the feasible range becomes either $(2\alpha_k^*, 0)$ or $(0, 2\alpha_k^*)$, depending on the sign of α_k^* . The negative sign is not too surprising, since if the function is locally concave, we do need an ascent direction to pass the hill, otherwise stuck at the local minimum. This aspect shares a similar flavour to the so-called 'gradient descent ascent' method for solving min-max problems, see e.g. Lin et al. (2020); Zheng et al. (2024).

Despite our non-convex applicability, the situation is highly challenging. Unlike the convex case where α_k^* is lower bounded, here it can take an arbitrary value. The worst case is when $\alpha_k^* = 0$, implying an empty selection range. This arises when $x^* - x^k \perp \nabla f(x^k)$, and a stepsize that can improve the current iterate x^k does not exist. On the other hand, if one can exclude such an orthogonal case, then convergence to the global optimal point is guaranteed, see Theorem 2.1.

While our theory is powerful, it is not instantly useful in practice, due to quantity $\langle x^* - x^k, -\nabla f(x^k) \rangle$ is not a priori knowledge. Experts may instantly realize that, by Taylor expansion, it is an upper bound for $f(x^k) - f(x^*)$ in the convex case, and the only concern is regarding $f(x^*)$. We show that,

(i) when $f(x^*) = 0$, the simplest tune-free stepsize $f(x^k)/||\nabla f(x^k)||^2$ is applicable. It is at least 1/(2L) large, see Proposition 3.2. In a special case, its two-times scaled version is optimal, see Section 3.1.3.

(ii) when $f(x^*)$ not known in advance, a parameter \overline{f}_0 is introduced as an initial guess for $f(x^*)$. It will be updated if some criteria violated, see details in Algorithm 1. Moreover, such a guess can be easily picked, for example, let $\overline{f}_0 = 0.1 \cdot f(x^0)$, where $f(x^0)$ is the initial objective value.

An outstanding benefit of our scheme is regarding the ill-conditioning issue, which is a well-known bottleneck for the GD algorithm. This aspect has been nicely illustrated in (Boyd & Vandenberghe, 2004, Sec. 9.3.2) through an \mathbb{R}^2 example, where an exact linear rate with factor $(\gamma - 1)^2/(\gamma + 1)^2$ is given, using an exact line search stepsize. A large γ (ill-conditioning) causes such a factor close to 1, implying the error has almost no change as GD iterating. Ours yields a factor of $(\gamma - 1)^2/(2\gamma^2 + 2)$, which is strictly smaller than 1/2, i.e., the error is at least halved each iteration, see more details in Section 4.1.

For notations, $\|\cdot\|$ denotes the Euclidean norm, induced by the inner product $\langle\cdot,\cdot\rangle$. The uppercase bold, lowercase bold, and not bold letters are used for matrices, vectors, and scalars, respectively.

094 095 096

1.1 LITERATURE: ADAPTIVE STEPSIZE

Here, we briefly discuss some developments of the stepsize adaption technique in the machine 097 learning field. The most popular family includes AdaGrad Duchi et al. (2011), RMSProp Tieleman & 098 Hinton. (2012), and Adam Kingma & Ba (2015). These approaches are strongly related to each other 099 and are heuristic methods that typically require tuning multiple parameters. Recently, Baydin et al. 100 (2018) propose to adaptively update the stepsize via a so-called 'hyper-gradient', which computes a 101 derivative over the stepsize parameter. The good news is that doing so adds very limited cost owing 102 to an element-wise product. The bad news is that the 'hyper-gradient' introduces a 'hyper-stepsize' 103 which still needs tuning (but tends to be easier). Also, a theoretical convergence guarantee is not yet 104 available. A follow-up work by Chandra et al. (2022) addresses the tuning issue by computing an 105 additional 'hyper-gradient' on the original 'hyper-stepsize'. This would introduce another 'hyperstepsize', and they apply the same procedure again, and so on, ad infinitum. The good news is that 106 each additionally introduced 'hyper-gradient' reduces the stepsize sensitivity, and eventually they can 107 easily pick an initial hyper-stepsize.

In view of these methods, we note that there is always an initial stepsize tuning issue, also referred to as 'the global learning rate' selection. This issue is avoided in our approach, since all of our choices, including the initial one, are mathematically computed.

1.2 KEY RESULTS

Below, we summarize 3 versions of our adaptive stepsize choices.

• (i) Theoretically, the k-th optimal choice

$$\alpha_k^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2} = \frac{\|\boldsymbol{x}^{\star} - \boldsymbol{x}^k\|}{\|\nabla f(\boldsymbol{x}^k)\|} \cos \eta_k, \quad k = 0, 1, \dots,$$
(1.1)

where $\eta_k = \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^k\| \| \nabla f(\boldsymbol{x}^k) \|}$. It admits an exact linear rate, with or without convexity: $\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^*\|^2 = (\Pi_{t=0}^k \sin^2 n_t) \|\boldsymbol{x}^0 - \boldsymbol{x}^*\|^2$

$$\boldsymbol{x}^{\kappa+1} - \boldsymbol{x}^{\star} \|^{2} = \left(\prod_{t=0}^{\kappa} \sin^{2} \eta_{t} \right) \| \boldsymbol{x}^{0} - \boldsymbol{x}^{\star} \|^{2}.$$
(1.2)

• (ii) The k-th practical-use choice (general version)

$${}^{\dagger}_{k} = \gamma_{0} \cdot \frac{f(\boldsymbol{x}^{k}) - \bar{f}_{0}}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}},$$
(1.3)

which admits the following exact linear rate, with or without convexity:

 α

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} = (\Pi_{t=0}^{k} \ \delta_{t}) \|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}, \qquad k = 0, 1, \dots,$$
(1.4)

where

$$\delta_t = 1 - \frac{\gamma_0}{\sigma_t} \left(2 - \frac{\gamma_0}{\sigma_t} \right) \cos^2 \eta_t, \qquad \sigma_t = \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{f(\boldsymbol{x}^t) - \bar{f}_0}.$$
 (1.5)

• (iii) The simplest practical-use choice (tune-free)

$$\widetilde{\alpha}_k = \frac{f(\boldsymbol{x}^k)}{\|\nabla f(\boldsymbol{x}^k)\|^2},\tag{1.6}$$

which guarantees convergence if f is convex and $f(x^*) = 0$. Empirically, it also works nicely for the non-convex MNIST problem in some settings.

ADAPTIVE STEPSIZE THEORY

Consider the following problem:

$$\underset{\boldsymbol{x} \in \mathbb{R}^n}{\operatorname{minimize}} f(\boldsymbol{x}), \tag{2.1}$$

where function $f: \mathbb{R}^n \to \mathbb{R}$ is assumed to be everywhere differentiable. The associated gradient descent (GD) iterates are

$$x^{k+1} = x^k - \alpha_k \nabla f(x^k), \quad k = 0, 1, \dots$$
 (2.2)

Throughout the rest of the paper, we assume $\nabla f(x^k) \neq 0$, unless GD already converged $x^k = x^*$. This assumption is necessary, since otherwise GD yields $x^{k+1} = x^k - \alpha_k \cdot \mathbf{0} = x^k$, and the stepsize selection issue becomes trivial.

2.1 SELECTION RANGE

First, we show a feasible selection range for stepsize α to guarantee convergence.

Proposition 2.1 (range). Consider GD in equation 2.2. While iterates not converged, let stepsize

$$\alpha_{k} \in \left(\frac{2\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}, 0\right) \bigcup \left(0, \frac{2\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}\right), \quad k = 0, 1, \dots$$
(2.3)

If such α_k exists $\forall k$. Then, convergence to the global optimal point is guaranteed.

Corollary 2.1. α_k as in equation 2.3 does not exist if and only if

$$\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle = 0.$$
 (2.4)

Remarks 2.1 (interpretation). In view of Corollary 2.1, it says that a feasible stepsize does not exist, if vectors $x^{\star} - x^{k}$ and $-\nabla f(x^{k})$ are orthogonal (zero vector case omitted by assumption). This is not surprising, since when orthogonality arises, by changing stepsize α_k alone, the future iterate $x^{k+1} = x^k - \alpha_k \nabla f(x^k)$ cannot be any closer to x^* than that of x^k .

162 2.2 OPTIMAL CHOICE

166

167 168 169

172

175 176 177

178

183 184

188 189

190

193 194

195 196

197

198 199

Here, we present the optimal stepsize choice from the above feasible range. It turns out to be its central point.

Theorem 2.1 (optimal choice). Consider GD in equation 2.2. The optimal k-th choice is given by

$$\alpha_k^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2} = \frac{\|\boldsymbol{x}^{\star} - \boldsymbol{x}^k\|}{\|\nabla f(\boldsymbol{x}^k)\|} \cos \eta_k,$$
(2.5)

170 where $\eta_k \stackrel{\text{\tiny def}}{=} \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^k\| \| \nabla f(\boldsymbol{x}^k) \|}$. It admits the following exact adaptive linear rate: 171 $\|\boldsymbol{x}^k + \boldsymbol{x}^k\| \| \nabla f(\boldsymbol{x}^k) \|$

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} = \left(\Pi_{t=0}^{k} \sin^{2} \eta_{t}\right) \|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}, \quad k = 0, 1, \dots$$
(2.6)

173 Remarks 2.2 (scaling invariance). GD equipped with α_k^* in equation 2.5 is invariant under a linearly 174 transformed function, $g(\cdot) = \rho f(\cdot), \forall \rho \neq 0$, since

$$\boldsymbol{x}^{k+1} = \boldsymbol{x}^{k} - \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\rho \nabla f(\boldsymbol{x}^{k}) \rangle}{\|\rho \nabla f(\boldsymbol{x}^{k})\|^{2}} \rho \nabla f(\boldsymbol{x}^{k}) = \boldsymbol{x}^{k} - \alpha_{k}^{\star} \nabla f(\boldsymbol{x}^{k}).$$
(2.7)

2.3 CONVEXITY

Suppose function f is convex. Then, much stronger guarantees and simplifications are available.
Corollary 2.2. Consider GD in equation 2.2. Suppose function f is convex. While iterates not converged, let stepsize

$$\alpha_k \in \left(0, \frac{2\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2}\right) = (0, 2\alpha_k^{\star}), \quad k = 0, 1, \dots$$
(2.8)

185 Then, the GD iterations are guaranteed to converge to the optimal point. 186 Remarks 2.3. Given a convex function f, relation $\langle x^* - x^k, -\nabla f(x^k) \rangle > 0$ always holds, unless 187 $x^k = x^*$.

2.3.1 L-SMOOTH

Here, we provide some characterizations via the *L*-smoothness assumption.

Definition 2.1. A differentiable convex function $f : \mathbb{R}^n \to \mathbb{R}$ is said to be L-smooth if

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\| \le L \|\boldsymbol{x} - \boldsymbol{y}\|, \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n.$$
 (2.9)

Proposition 2.2. Suppose function $f : \mathbb{R}^n \to \mathbb{R}$ is L-smooth. Then,

$$\alpha_k^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2} \ge \frac{1}{L}, \quad k = 0, 1....$$
(2.10)

Corollary 2.3. The fixed stepsize selection range is a subset of our adaptive one, i.e.,

$$\left(0, \frac{2}{L}\right) \subseteq \left(0, 2\alpha_k^{\star}\right), \quad k = 0, 1....$$

$$(2.11)$$

200 201 202

203

206

207 208

214 215

3 PRACTICAL USE

The above theory involves optimal point x^* , hence not instantly useful in practice. Here, we address it via approximation.

Theorem 3.1. Consider GD in equation 2.2. While iterates not converged, we propose stepsize

$$\alpha_k^{\dagger} = \gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - f_0}{\|\nabla f(\boldsymbol{x}^k)\|^2},$$
(3.1)

where γ_0 is a tunable parameter; \overline{f}_0 is a guessed smallest objective value. It admits the following exact linear rate:

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} = \left(\prod_{t=0}^{k} \delta_{t}\right) \|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2},$$
(3.2)

212 213 where

$$\delta_t = 1 - \frac{\gamma_0}{\sigma_t} \left(2 - \frac{\gamma_0}{\sigma_t} \right) \cos^2 \eta_t, \qquad \sigma_t = \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{f(\boldsymbol{x}^t) - \bar{f}_0}, \tag{3.3}$$

and where
$$\eta_t = \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^t\| \|\nabla f(\boldsymbol{x}^t)\|}$$

Corollary 3.1 (convergence). *While GD iterates not converged, let*

$$\gamma_0 \in (2\sigma_k, 0) \cup (0, 2\sigma_k), \quad k = 0, 1, \dots$$
 (3.4)

If such γ_0 exists $\forall k$, then

$$\delta_k = 1 - \frac{\gamma_0}{\sigma_k} \left(2 - \frac{\gamma_0}{\sigma_k} \right) \cos^2 \eta_k \in (0, 1), \quad \forall k,$$
(3.5)

which guarantees convergence to the global optimal point.

Corollary 3.2. The optimal k-th choice of the tunable parameter γ_0 is

$$\gamma_0^{\star} = \underset{\gamma_0}{\operatorname{argmax}} \frac{\gamma_0}{\sigma_k} \left(2 - \frac{\gamma_0}{\sigma_k} \right) = \sigma_k.$$
(3.6)

In this case, the rate factor

$$\delta_k^{\star} = 1 - \frac{\gamma_0^{\star}}{\sigma_k} \left(2 - \frac{\gamma_0^{\star}}{\sigma_k} \right) \cos^2 \eta_k = \sin^2 \eta_k, \tag{3.7}$$

implying optimality attained (recall Theorem 2.1), i.e., exact approximation.
 230

Remarks 3.1. In view of Corollary 3.2, the approximation is exact if one can adaptively select $\gamma_0 = \sigma_k$, $\forall k$. There does exist a special case where σ_k is a known constant, see Section 3.1.3. However, in general, we do not know σ_k in advance, and our approximation hence not exact. Also, for ease of use, we typically fix γ_0 to be a constant, which is theoretically sub-optimal.

3.1 CONVEXITY

Suppose function f is convex. Then, we have stronger guarantees and a tune-free stepsize selection scheme.

Corollary 3.3 (convergence). Suppose function f is convex. While GD iterates not converged, let $\gamma_0 \in (0, 2\sigma_k), \quad \forall k,$ (3.8)

241
242 where
$$\sigma_k = \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{f(\boldsymbol{x}^k) - \bar{f}_0}$$
. Then, the rate factor satisfies

243 244

245

246

248

252 253

255 256

257 258

265 266

218

219

220

221

222

227 228

231

232

233

234 235

236

$$\delta_k = 1 - \frac{\gamma_0}{\sigma_k} \left(2 - \frac{\gamma_0}{\sigma_k} \right) \cos^2 \eta_k \in (0, 1), \quad \forall k,$$
(3.9)

which guarantees convergence.

247 3.1.1 TUNE-FREE CASE

Here, we require full knowledge of $f(x^*)$.

Proposition 3.1. Consider GD in equation 2.2. Suppose function f is convex, with optimal objective value $f(x^*)$ known in advance. Then, stepsize

$$\widetilde{\alpha}_k = \gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - f(\boldsymbol{x}^\star)}{\|\nabla f(\boldsymbol{x}^k)\|^2}, \quad \gamma_0 \in (0, 2],$$
(3.10)

254 guarantees convergence, with an exact linear rate:

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} = \left(\Pi_{t=0}^{k} \ \delta_{t}\right) \|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}, \qquad (3.11)$$

where

$$\delta_t = 1 - \frac{\gamma_0}{\sigma_t} \left(2 - \frac{\gamma_0}{\sigma_t} \right) \cos^2 \eta_t, \quad \sigma_t = \frac{\langle \boldsymbol{x}^\star - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{f(\boldsymbol{x}^t) - f(\boldsymbol{x}^\star)}, \tag{3.12}$$

259 260 where $\eta_t = \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^t\| \|\nabla f(\boldsymbol{x}^t)\|}$

Remarks 3.2. The above tune-free case can happen in practice. A typical example is when $f(x^*) = 0$, arising in (i) solving a huge-scale linear system Ax = b, where A^{-1} is too expensive to calculate directly; (ii) f is a loss function with zero-loss at the optimal point, as in many classification problems. **Corollary 3.4.** Suppose f is a non-linear convex function. Then, when $x^k \neq x^*$, we have

$$\gamma_0^{\star} = \sigma_k = \frac{\left\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle}{f(\boldsymbol{x}^k) - f(\boldsymbol{x}^{\star})} > 1, \quad \forall k.$$
(3.13)

267 *Remarks* 3.3. equation 3.13 follows instantly from Taylor expansion. It implies that we should choose 268 $\gamma_0 > 1$ in our convex tune-free case. However, it does not tell exactly how much larger than 1, our 269 default choice is therefore conservatively set to $\gamma_0 = 1$. Additionally, we assume *f* being non-linear, since minimizing a linear or affine function is trivial (unbounded below).

270 3.1.2 *L*-SMOOTH

272 Here, we provide some characterizations via the *L*-smooth assumption.

Proposition 3.2. Suppose function $f : \mathbb{R}^n \to \mathbb{R}$ is L-smooth. Then,

$$\frac{f(\boldsymbol{x}^k) - f(\boldsymbol{x}^\star)}{\|\nabla f(\boldsymbol{x}^k)\|^2} \ge \frac{1}{2L},\tag{3.14}$$

Proposition 3.3 (optimality gap). Let function $f : \mathbb{R}^n \to \mathbb{R}$ be L-smooth. Then,

$$\underbrace{\frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}}_{optimal} - \underbrace{\frac{f(\boldsymbol{x}^{k}) - f(\boldsymbol{x}^{\star})}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}}_{estimated(\gamma_{0}=1)} \geq \frac{1}{2L}.$$
(3.15)

Remarks 3.4. The positive gap from Proposition 3.3 with $\gamma_0 = 1$ is not surprising, since we already seen from Corollary 3.4 that the optimal parameter γ_0^* is strictly larger than 1 (and γ_0^* attains optimality by Corollary 3.2). The result here is strengthened, with the gap characterized by *L*, instead of only being positive.

287 3.1.3 PRACTICAL EXACT APPROXIMATION

Here, we show special cases that our practical-use stepsize choice attains the theoretical optimum, by simply selecting $\gamma_0 = 2$. Consider

$$\underset{\boldsymbol{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|^2.$$
(3.16)

where $oldsymbol{x} \in \mathbb{R}^n, oldsymbol{b} \in \mathbb{R}^m, oldsymbol{A} \in \mathbb{R}^{m imes n}.$

(i) Suppose A is a full-rank square matrix. We have $x^* = A^{-1}b$. It follows that,

$$\alpha_{k}^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}} = \frac{\langle \boldsymbol{A}^{-1}\boldsymbol{b} - \boldsymbol{x}^{k}, -\boldsymbol{A}^{T}(\boldsymbol{A}\boldsymbol{x}^{k} - \boldsymbol{b}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}} = \frac{\|\boldsymbol{A}\boldsymbol{x}^{k} - \boldsymbol{b}\|^{2}}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}} = \frac{2 \cdot f(\boldsymbol{x}^{k})}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}},$$
(3.17)

corresponding to our practical-use stepsize with $\gamma_0 = 2$ and $f(\mathbf{x}^*) = 0$, recall equation 3.10.

(ii) Suppose b = 0. We have $x^* = 0$. It follows that,

$$\alpha_{k}^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}} = \frac{\langle \boldsymbol{0} - \boldsymbol{x}^{k}, -\boldsymbol{A}^{T}(\boldsymbol{A}\boldsymbol{x}^{k} - \boldsymbol{0}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}} = \frac{\|\boldsymbol{A}\boldsymbol{x}^{k}\|^{2}}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}} = \frac{2 \cdot f(\boldsymbol{x}^{k})}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}},$$
(3.18)

which is similar to the above case.

3.2 GENERAL PRACTICAL USE ALGORITHM

Here, we consider f_0 being an inaccurate guess. It will be updated if certain criteria violated.

Algorithm 1 Linear-rate gradient decent (auto correction version)

```
312
                  Input: initialization x^0; iteration number counter k = 0;
313
                  Input: guessed \bar{f}_0, tunable parameter \gamma_0;
314
                  Input: shrinking factors \tau_1, \tau_2 \in (0, 1), threshold T.
315
                    1: while iterates not converged do
316
                    2:
                                    k \leftarrow k+1
                               \begin{array}{cccc} & & & & & & & \\ \alpha_k & \leftarrow & & & \\ \alpha_k & \leftarrow & & & \\ \gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{\|\nabla f(\boldsymbol{x}^k)\|^2}, \\ \boldsymbol{x}^{k+1} & \leftarrow & & & \\ \boldsymbol{x}^k - & \alpha_k \nabla f(\boldsymbol{x}^k) \end{array}
317
                    3:
318
                    4:
319
                               Correction:
                    5:
320
                               If f(\boldsymbol{x}^{k+1}) > T \cdot f(\boldsymbol{x}^k), set \gamma_0 \leftarrow \tau_1 \cdot \gamma_0 and \boldsymbol{x}^{k+1} \leftarrow \boldsymbol{x}^k.
321
                               If \alpha_k \leq 0, set \bar{f}_0 \leftarrow \tau_2 \cdot \bar{f}_0.
322
                    6: end while
323
                  Output: x^{k+1}
```

4 EXAMPLES

4.1 \mathbb{R}^2 quadratic program

Here, we consider a simple example from (Boyd & Vandenberghe, 2004, Sec. 9.3.2):

$$\underset{x_{1},x_{2}}{\text{minimize}} \ \frac{1}{2} \left(x_{1}^{2} + \gamma x_{2}^{2} \right), \tag{4.1}$$

where $\gamma > 0$, $\boldsymbol{x} = [x_1, x_2]^T$. We employ initialization $\boldsymbol{x}^0 = [\gamma, 1]^T$. For this problem, the Lipschitz constant L and the condition number both equal to γ , and the conditioning state is fully tractable.

In this section, through some specific examples, we demonstrate the power of our adaptive stepsize.

Below, we compare our approach with the exact line search method, which finds a stepsize choice via

$$\alpha_k = \operatorname*{argmin}_{\alpha_k > 0} f(\boldsymbol{x}^k - \alpha_k \nabla f(\boldsymbol{x}^k)). \tag{4.2}$$

• Following from (Boyd & Vandenberghe, 2004, Sec. 9.3.2), the *k*-th iterate with an exact line search stepsize is given by

$$k_1^k = \gamma \left(\frac{\gamma - 1}{\gamma + 1}\right)^k, \quad x_2^k = \left(-\frac{\gamma - 1}{\gamma + 1}\right)^k,$$
(4.3)

with an exact convergence rate

x

$$\frac{\|\boldsymbol{x}^{k} - \boldsymbol{x}^{\star}\|^{2}}{\|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}} = \left(\frac{\gamma - 1}{\gamma + 1}\right)^{2k}.$$
(4.4)

If γ is large (ill-conditioning), the above factor is close to 1, i.e., $\|\boldsymbol{x}^k - \boldsymbol{x}^\star\|^2$ is similar to $\|\boldsymbol{x}^0 - \boldsymbol{x}^\star\|^2$. That said, the ground-truth error has little changes after k iterations.

• Our optimal choice α_k^{\star} yields

$$x_{1}^{k} = \gamma^{k_{2}-k_{1}+1} \left(\frac{\gamma-1}{2}\right)^{k_{1}} \left(\frac{\gamma-1}{\gamma^{2}+1}\right)^{k_{2}},$$

$$x_{2}^{k} = (-1)^{k_{1}+k_{2}} \left(\frac{\gamma-1}{2}\right)^{k_{1}} \left(\frac{\gamma-1}{\gamma^{2}+1}\right)^{k_{2}},$$
(4.5)

where $k_1 \stackrel{\text{def}}{=} \lfloor \frac{k+1}{2} \rfloor$, $k_2 \stackrel{\text{def}}{=} \lfloor \frac{k}{2} \rfloor$, and where $\lfloor \cdot \rfloor$ denotes the floor operation (the closest smaller integer). Ours admits the following convergence rate factor:

$$\frac{\|\boldsymbol{x}^{k} - \boldsymbol{x}^{\star}\|^{2}}{\|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}} = \frac{\gamma^{2(k_{2} - k_{1} + 1)} + 1}{\gamma^{2} + 1} \left(\frac{\gamma - 1}{2}\right)^{2k_{1}} \left(\frac{\gamma - 1}{\gamma^{2} + 1}\right)^{2k_{2}} = \left(\frac{1}{2}\right)^{k} \left(\frac{\gamma - 1}{\sqrt{\gamma^{2} + 1}}\right)^{2k}.$$
 (4.6)

Since $\gamma > 0$, we have $(\gamma - 1)/\sqrt{\gamma^2 + 1} < 1$. Our factor is therefore strictly smaller than $1/2^k$.



Figure 1: exact line search vs. our stepsize, with conditioning controlled by γ .

4.1.1 STRICT BETTER PERFORMANCE

Here, we show that our rate is strictly better than the one for exact line search, except when $\gamma = 1$ both methods converge in exactly one iteration.

To this end, suppose $\gamma \neq 1$. Divide our rate factor by that of the exact line search, we arrive at

$$\frac{\delta_{\text{ours}}^{(k)}}{\delta_{\text{line-search}}^{(k)}} = \left(\frac{1}{2}\right)^k \left(\frac{\gamma - 1}{\sqrt{\gamma^2 + 1}}\right)^{2k} \left(\frac{\gamma + 1}{\gamma - 1}\right)^{2k} = \left(\frac{\gamma^2 + 2\gamma + 1}{2\gamma^2 + 2}\right)^k < 1, \tag{4.7}$$

where the last inequality follows from the denominator being larger when $\gamma \neq 1$, since

$$2\gamma^{2} + 2 - (\gamma^{2} + 2\gamma + 1) = \gamma^{2} - 2\gamma + 1 = (\gamma - 1)^{2} > 0.$$
(4.8)

4.1.2 INSTANT CONVERGENCE

Here, we perform an additional test on our rate factor $\sin^2 \eta_k$. An observation is that if $\sin \eta_k = 0$, GD must converge instantly. In the current example, we can easily verify it using a sparse initialization, say $(x_1^0, x_2^0) = (50, 0)$. Indeed, 1-step instant convergence is observed.



Figure 2: zero angle case, $\gamma = 10$.

4.2 GEOMETRIC PROGRAM

Here, we consider an unconstrained geometric program from (Boyd & Vandenberghe, 2004, Sec. 9.3), and our Algorithm 1 will apply. Consider

$$\underset{\boldsymbol{x}}{\text{minimize }} \log \left(\sum_{i=1}^{m} \exp \left(\boldsymbol{a}_{i}^{T} \boldsymbol{T} \boldsymbol{x} + b_{i} \right) \right), \tag{4.9}$$

where $\boldsymbol{x} \in \mathbb{R}^n$, $\boldsymbol{a}_i \in \mathbb{R}^n$, $b_i \in \mathbb{R}$, and $\boldsymbol{T} = \operatorname{diag}\left(\left[1, \gamma^{\frac{1}{n}}, \gamma^{\frac{2}{n}}, ..., \gamma^{\frac{n-1}{n}}\right]\right)$ is a diagonal matrix that promotes ill-conditioning.

Below, we compare our Algorithm 1 with (i) a fine-tuned fixed stepsize; (ii) a fine-tuned Nesterov's accelerated gradient descent (N-AGD) Nesterov (1983). The tuning is performed on a fine grid with a fixed random number generator, hence shows roughly their best performances. Our parameters are very roughly picked as $\gamma_0 = 1, \tau_1 = \tau_2 = 0.5, T = 1, f_0 = 0.1 \cdot f(x^0)$ and no further tuning.



Remarks 4.1 (worst-case acceleration). We observe that N-AGD provides almost no acceleration in 433 the ill-conditioning setting here. Let us note that its well-known $O(1/k^2)$ rate is only guaranteed in 434 a worst-case sense and does not necessarily accelerate in practice, see a discussion in (Ryu & Yin, 435 2022, Sec. 12.3).

Remarks 4.2. Due to rough choice of parameters, our guessed $\bar{f}_0 = 0.1 \cdot f(x^0)$ admits a significant 437 performance gap compared to an ideal tune-free case $f_0 = f(x^*)$ (not a priori knowledge). How to 438 improve such a gap is left for future work.

440 4.3 NON-CONVEX MNIST

Here, we consider the MNIST classification problem via a 2-layer neural network, with ReLu activation, 200 hidden units, and softmax loss function. Following the literature, we consider a mini-batch setting. We compare ours with the state-of-the-art algorithms, Nesterov's accelerated gradient descent (N-AGD) Nesterov (1983) and Adaptive moment estimation (Adam) Kingma & Ba (2015).

4.3.1 TUNE-FREE CASE

We start with a special case that stepsize $\alpha_k = f(\mathbf{x}^k) / \|\nabla f(\mathbf{x}^k)\|^2$ alone works nicely. We consider minimizing the softmax loss only (no regularization) under a relatively large mini-batch size.



Figure 4: Our tune-free case, with mini-batch size 1024.

Remarks 4.3. Fig 4a and Fig 4b record the training and validation accuracies, respectively. We observe that they share a highly similar trend (but not the same). Ours exhibits consistent advantages over the others.

4.3.2 GENERAL CASE

Here, we consider a general case, minimizing softmax loss function with l_2 -norm regularization (on the weights). We adopt a commonly used mini-batch size of 128. Our Algorithm 1 is applied, with roughly picked parameters $\bar{f}_0 = 0$, $\gamma_0 = 1$, T = 5, $\tau_1 = 0.25$ (τ_2 omitted).





Figure 5: General case with l_2 -norm regularization; mini-batch size 128.

Remarks 4.5. Ours only has advantage in the validation stage, where consistently higher accuracy is observed. Luckily, the validation accuracy is all we need, hence ours remains a better choice. Additionally, we suspect our advantage can be enlarged if more careful parameter choices are employed, which is left for future research.

490 491 492

507

508 509

510

511

522

5 CONCLUSION

493 In this work, we established a general theory on the adaptive stepsize selection issue, including 494 feasible selection range, convergence rate, and optimal choice. Specifically, in the convex case, we 495 show an adaptive range $(0, 2\alpha_k^*)$ that guarantees convergence, which enlarges the classical fixed 496 one (0, 2/L). Its centre α_k^{\star} is the optimal choice, admitting an exact linear rate with factor $\sin^2 \eta_k$. 497 Our theory also applies to a non-convex function, except the situation is much more challenging. 498 The optimal stepsize can now be negative, and the feasible range set could be empty when some 499 orthogonality arises. On the other hand, if a feasible stepsize choice always exists, then convergence 500 to the global optimal point is guaranteed.

Despite the great power of our theory, it involves some optimal point information. To enable its
 practical use, we propose an approximation strategy. Such an approximation can be exact in a special
 practical scenario but in general sub-optimal. It also admits an exact linear convergence rate, and we
 numerically test its power through several examples. Outstandingly, a tune-free version works nicely
 for the non-convex MNIST problem via neural networks.

6 Reproducibility Statement

All figures in this manuscript can be reproduced via the MATLAB codes submitted as supplementary material.

512 513 REFERENCES

- Luís B. Almeida, Thibault Langlois, José D. Amaral, and Alexander Plakhov. *Parameter adaptation in stochastic optimization*, pp. 111–134. Cambridge University Press, USA, 1999. ISBN 0521652634.
- 517 Cem Anil, James Lucas, and Roger Grosse. Sorting out Lipschitz function approximation. In
 518 Proceedings of the 36th International Conference on Machine Learning, pp. 291–301, 2019.
- Atilim Gunes Baydin, Robert Cornish, David Martínez Rubio, Mark Schmidt, and Frank D. Wood.
 Online learning rate adaptation with hypergradient descent. In *Sixth International Conference on Learning Representations, ICLR*, 2018.
- Yoshua Bengio. Practical recommendations for gradient-based training of deep architectures. In
 Neural Networks: Tricks of the Trade: Second Edition, pp. 437–478. Springer, 2012.
- Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Kartik Chandra, Audrey Xie, Jonathan Ragan-Kelley, and ERIK MEIJER. Gradient descent: The ultimate optimizer. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 8214–8225. Curran Associates, Inc., 2022.
- Patrick L. Combettes and Jean-Christophe Pesquet. Lipschitz certificates for layered network
 structures driven by averaged activation operators, 2020.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and
 stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011.
- Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas. Efficient and accurate estimation of lipschitz constants for deep neural networks. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
 - Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep Learning. MIT Press, 2016.

540 541 542	Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, 2015.
543 544 545 546	Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In Hal Daumé III and Aarti Singh (eds.), <i>Proceedings of the 37th International Conference on Machine Learning</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pp. 6083–6093. PMLR, 13–18 Jul 2020.
547 548	Yurii Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. Proceedings of the USSR Academy of Sciences, 269:543–547, 1983.
549 550 551	Yurii Nesterov. <i>Lectures on Convex Optimization</i> . Springer Publishing Company, Incorporated, 2nd edition, 2018. ISBN 3319915770.
552 553	Sebastian Ruder. An overview of gradient descent optimization algorithms. <i>ArXiv</i> , abs/1609.04747, 2016.
555 556	David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. In <i>Nature</i> , 1986.
557 558	Ernest K Ryu and Wotao Yin. Large-scale convex optimization: algorithms & analyses via monotone operators. Cambridge University Press, 2022.
559 560 561	T. Tieleman and G. Hinton. Lecture 6.5 – RMSProp: Divide the gradient by a running average of its recent magnitude. In <i>COURSERA: Neural Networks for Machine Learning</i> , 2012.
502 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 584 585 586 587 588 589	Taoli Zheng, Linglingzhi Zhu, Anthony Man-Cho So, José Blanchet, and Jiajin Li. Universal gradient descent ascent method for nonconvex-nonconcave minimax optimization. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems</i> , 2024.
591 592	

Appendix А

The gradient descent (GD) iterates are

$$\boldsymbol{x}^{k+1} = \boldsymbol{x}^k - \alpha_k \nabla f(\boldsymbol{x}^k), \quad k = 0, 1, \dots$$
(A.1)

We assume $\nabla f(x^k) \neq 0$, unless $x^k = x^*$. This assumption is necessary, since otherwise stepsize selection becomes trivial.

A.1 PROOF OF PROPOSITION 2.1

Our Proposition 2.1, restated here as

Proposition A.1 (range). Consider GD in equation A.1. While iterates not converged, let stepsize

$$\alpha_{k} \in \left(\frac{2\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}, 0\right) \bigcup \left(0, \frac{2\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}\right), \quad k = 0, 1, \dots \quad (A.2)$$

If such α_k exists $\forall k$. Then, convergence to the global optimal point is guaranteed.

X

Proof. Let us note that

$$\begin{aligned} & \| \boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star} \|^{2} - \| \boldsymbol{x}^{k} - \boldsymbol{x}^{\star} \|^{2} &= - \| \boldsymbol{x}^{k+1} - \boldsymbol{x}^{k} \|^{2} - 2 \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k+1}, \boldsymbol{x}^{k+1} - \boldsymbol{x}^{k} \rangle, \\ & = - (\alpha_{k})^{2} \| \nabla f(\boldsymbol{x}^{k}) \|^{2} - 2 \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k} + \alpha_{k} \nabla f(\boldsymbol{x}^{k}), -\alpha_{k} \nabla f(\boldsymbol{x}^{k}) \rangle, \\ & = \alpha_{k}^{2} \| \nabla f(\boldsymbol{x}^{k}) \|^{2} - 2 \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\alpha_{k} \nabla f(\boldsymbol{x}^{k}) \rangle, \\ & = \alpha_{k} \left(\alpha_{k} \| \nabla f(\boldsymbol{x}^{k}) \|^{2} + 2 \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle \right). \end{aligned}$$
(A.3)

All we need is the above right-hand side being negative, implying the ground-truth error is strictly decreasing, hence guarantees convergence. This yields equation A.2, which involves one empty set, depends on the sign of the term $\langle x^* - x^k, -\nabla f(x^k) \rangle$. The proof is now concluded.

A.2 PROOF OF THEOREM 2.1

Our Theorem 2.1, restated here as

Theorem A.1 (optimal choice). Consider GD in equation A.1. The optimal k-th choice is given by

$$\alpha_k^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2} = \frac{\|\boldsymbol{x}^{\star} - \boldsymbol{x}^k\|}{\|\nabla f(\boldsymbol{x}^k)\|} \cos \eta_k, \tag{A.4}$$

where $\eta_k = \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^k\| \|\nabla f(\boldsymbol{x}^k)\|}$. It admits the following exact adaptive linear rate:

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} = \left(\Pi_{t=0}^{k} \sin^{2} \eta_{t}\right) \|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}, \quad k = 0, 1, \dots$$
(A.5)

Proof. Following from equation A.3, we would like its right-hand side term as negative as possible, which leads to

$$\underset{\alpha_k}{\text{minimize }} (\alpha_k)^2 \left\| \nabla f(\boldsymbol{x}^k) \right\|^2 - 2\alpha_k \left\langle \boldsymbol{x}^* - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle.$$
(A.6)

Its solution is

$$\alpha_k^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2} = \frac{\|\boldsymbol{x}^{\star} - \boldsymbol{x}^k\|}{\|\nabla f(\boldsymbol{x}^k)\|} \cos \eta_k, \tag{A.7}$$

where $\eta_k = \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^k\| \| \nabla f(\boldsymbol{x}^k) \|}$. Substituting it back to equation A.6, we obtain the minimal objective value being

$$(\alpha_k^{\star})^2 \left\| \nabla f(\boldsymbol{x}^k) \right\|^2 - 2\alpha_k^{\star} \left\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle = - \|\boldsymbol{x}^{\star} - \boldsymbol{x}^k\|^2 \cos^2 \eta_k.$$
(A.8)

At last, by equation A.3, we obtain

The proof is concluded by considering all iterations, from 0 to the current k-th one.

A.3 PROOF OF PROPOSITION 2.2

Lemma A.1 (Bailton-Haddad Theorem). Let function $f : \mathbb{R}^n \to \mathbb{R}$ be L-smooth. The following holds:

$$\frac{1}{L} \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\|^2 \le \langle \boldsymbol{x} - \boldsymbol{y}, \nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y}) \rangle, \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n.$$
(A.10)

Our Proposition 2.2, restated here as

Proposition A.2. Suppose function $f : \mathbb{R}^n \to \mathbb{R}$ is L-smooth. Then,

$$\alpha_k^{\star} = \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2} \ge \frac{1}{L}, \quad k = 0, 1....$$
(A.11)

Proof. By Lemma A.1, we have

$$\frac{1}{L} \|\nabla f(\boldsymbol{x}^{\star}) - \nabla f(\boldsymbol{x}^{k})\|^{2} \leq \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, \nabla f(\boldsymbol{x}^{\star}) - \nabla f(\boldsymbol{x}^{k}) \rangle.$$
(A.12)

Rearranging the terms concludes the proof.

A.4 PROOF OF THEOREM 3.1

Our Theorem 3.1, restated here as

Theorem A.2. Consider GD in equation A.1. While iterates not converged, we propose stepsize

$$\alpha_k^{\dagger} = \gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{\|\nabla f(\boldsymbol{x}^k)\|^2}, \tag{A.13}$$

where γ_0 is a tunable parameter; \overline{f}_0 is a guessed smallest objective value. It admits the following *exact linear rate:*

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} = \left(\Pi_{t=0}^{k} \ \delta_{t}\right) \|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}, \tag{A.14}$$

where

$$\delta_t = 1 - \frac{\gamma_0}{\sigma_t} \left(2 - \frac{\gamma_0}{\sigma_t} \right) \cos^2 \eta_t, \qquad \sigma_t = \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{f(\boldsymbol{x}^t) - \bar{f}_0}, \tag{A.15}$$

and where $\eta_t = \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^t\| \| \nabla f(\boldsymbol{x}^t) \|}.$

Proof. Recall error characterization from equation A.3

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} - \|\boldsymbol{x}^{k} - \boldsymbol{x}^{\star}\|^{2} = (\alpha_{k})^{2} \|\nabla f(\boldsymbol{x}^{k})\|^{2} - 2\alpha_{k} \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle.$$
(A.16)

Substituting α_k^{\dagger} in equation A.13 to the right-hand side above, yields

$$\begin{aligned} \text{r.h.s.} &= \left(\gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{\|\nabla f(\boldsymbol{x}^k)\|^2}\right)^2 \left\|\nabla f(\boldsymbol{x}^k)\right\|^2 - 2\gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{\|\nabla f(\boldsymbol{x}^k)\|^2} \left\langle \boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle, \\ &= \left(\gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{\|\nabla f(\boldsymbol{x}^k)\|}\right)^2 - 2\gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{\|\nabla f(\boldsymbol{x}^k)\|^2} \left\langle \boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle, \\ &= \left(\frac{\left\langle \boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle}{\|\nabla f(\boldsymbol{x}^k)\|}\right)^2 \left(\left(\gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{|\boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k)\rangle}\right)^2 - \frac{1}{2}\right)^2 \\ &= \left(\frac{\left\langle \boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle}{\|\nabla f(\boldsymbol{x}^k)\|}\right)^2 \left(\left(\gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{|\boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k)\rangle}\right)^2 - \frac{1}{2}\right)^2 \\ &= \left(\frac{\left\langle \boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle}{\|\nabla f(\boldsymbol{x}^k)\|}\right)^2 \left(\left(\gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - \bar{f}_0}{|\boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k)\rangle}\right)^2 - \frac{1}{2}\right)^2 \\ &= \left(\frac{\left\langle \boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \right\rangle}{\|\nabla f(\boldsymbol{x}^k)\|}\right)^2 \left(\frac{1}{2}\left(\frac{1}{2}\right)^2 + \frac{1}{2}\left(\frac{1}{2}\right)^2 + \frac{1}{$$

$$\left(\begin{array}{c} \|\nabla f(\boldsymbol{x}^{k})\| & f(\boldsymbol{x}^{k}) \| \\ 2\gamma_{0} \cdot \frac{f(\boldsymbol{x}^{k}) - \bar{f}_{0}}{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle} \right), \quad (A.17)$$

Let $\sigma_k = \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{f(\boldsymbol{x}^k) - \bar{f}_0}$. Invoke the l.h.s. of equation A.16, we arrive at $\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^2 - \|\boldsymbol{x}^k - \boldsymbol{x}^{\star}\|^2 = \left(\frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|}\right)^2 \left(\left(\frac{\gamma_0}{\sigma_k}\right)^2 - 2 \cdot \frac{\gamma_0}{\sigma_k}\right),$ $=\left(\frac{\langle \boldsymbol{x}^{\star}-\boldsymbol{x}^{k},-\nabla f(\boldsymbol{x}^{k})\rangle}{\|\nabla f(\boldsymbol{x}^{k})\|\|\boldsymbol{x}^{k}-\boldsymbol{x}^{\star}\|}\right)^{2}\left(\left(\frac{\gamma_{0}}{\sigma_{k}}\right)^{2}-2\cdot\frac{\gamma_{0}}{\sigma_{k}}\right)\|\boldsymbol{x}^{k}-\boldsymbol{x}^{\star}\|^{2},$ $= \cos^2 \eta_k \cdot \left(\left(\frac{\gamma_0}{\sigma_k} \right)^2 - 2 \cdot \frac{\gamma_0}{\sigma_k} \right) \| \boldsymbol{x}^k - \boldsymbol{x}^\star \|^2,$ (A.18) where $\eta_k = \arccos \frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\boldsymbol{x}^{\star} - \boldsymbol{x}^k\| \| \nabla f(\boldsymbol{x}^k) \|}$. It follows that $\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^2 = \left(1 - rac{\gamma_0}{\sigma_k} \left(2 - rac{\gamma_0}{\sigma_k}\right) \cos^2 \eta_k\right) \|\boldsymbol{x}^k - \boldsymbol{x}^{\star}\|^2.$ (A.19) Considering all iterations t = 0, 1, 2...k gives equation A.14. The proof is now concluded. A.5 PROOF OF PROPOSITION 3.1 Our Proposition 3.1, restated here as **Proposition A.3** (tune-free stepsize). Consider GD in equation 2.2. Suppose function f is convex, with optimal objective value $f(\mathbf{x}^*)$ known in advance. Then, any choice from below $\widetilde{\alpha}_k = \gamma_0 \cdot \frac{f(\boldsymbol{x}^k) - f(\boldsymbol{x}^\star)}{\|\nabla f(\boldsymbol{x}^k)\|^2}, \quad \gamma_0 \in (0, 2].$ (A.20) guarantees convergence, with an exact linear rate:

$$\|\boldsymbol{x}^{k+1} - \boldsymbol{x}^{\star}\|^{2} = \left(\Pi_{t=0}^{k} \ \delta_{t}\right) \|\boldsymbol{x}^{0} - \boldsymbol{x}^{\star}\|^{2}, \tag{A.21}$$

where

$$\delta_t = 1 - \frac{\gamma_0}{\sigma_t} \left(2 - \frac{\gamma_0}{\sigma_t} \right) \cos^2 \eta_t, \quad \sigma_t = \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{f(\boldsymbol{x}^t) - f(\boldsymbol{x}^*)}, \tag{A.22}$$

where $\eta_t = \arccos \frac{\langle \boldsymbol{x}^* - \boldsymbol{x}^t, -\nabla f(\boldsymbol{x}^t) \rangle}{\|\boldsymbol{x}^* - \boldsymbol{x}^t\| \|\nabla f(\boldsymbol{x}^t)\|}$.

Proof. Given $x^k \neq x^*$, the following holds:

$$f(\boldsymbol{x}^{\star}) > f(\boldsymbol{x}^{k}) + \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, \nabla f(\boldsymbol{x}^{k}) \rangle, \qquad (A.23)$$

where we exclude the case of a linear function, since it is unbounded below and is therefore trivial to minimize. Rearranging the terms, yields

$$f(\boldsymbol{x}^{k}) - f(\boldsymbol{x}^{\star}) < \left\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \right\rangle$$
(A.24)

Suppose $\gamma_0 \in (0, 2]$. Then,

$$y_0 \cdot \frac{f(\boldsymbol{x}^k) - f(\boldsymbol{x}^\star)}{\|\nabla f(\boldsymbol{x}^k)\|^2} < \frac{2\langle \boldsymbol{x}^\star - \boldsymbol{x}^k, -\nabla f(\boldsymbol{x}^k) \rangle}{\|\nabla f(\boldsymbol{x}^k)\|^2}.$$
(A.25)

It says that the left-hand side above always lies within the feasible range, recall equation A.2. Its convergence rate follows directly from Theorem A.2. The proof is now concluded.

A.6 PROOF OF PROPOSITION 3.2

Our Proposition 3.2, restated here as

Proposition A.4. Suppose function $f : \mathbb{R}^n \to \mathbb{R}$ is L-smooth. Then,

$$\frac{f(\boldsymbol{x}^k) - f(\boldsymbol{x}^\star)}{\|\nabla f(\boldsymbol{x}^k)\|^2} \ge \frac{1}{2L},\tag{A.26}$$

Proof. The *L*-smoothness assumption implies

$$f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \langle \boldsymbol{y} - \boldsymbol{x}, \nabla f(\boldsymbol{x}) \rangle + \frac{L}{2} \| \boldsymbol{y} - \boldsymbol{x} \|^2, \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n.$$
(A.27)

760 We may perform the following minimization:

minimize
$$f(\boldsymbol{x}) + \langle \boldsymbol{y} - \boldsymbol{x}, \nabla f(\boldsymbol{x}) \rangle + \frac{L}{2} \| \boldsymbol{y} - \boldsymbol{x} \|^2,$$
 (A.28)

and obtain a minimizer

$$\widehat{\boldsymbol{y}} = \boldsymbol{x} - \frac{1}{L} \nabla f(\boldsymbol{x}) \tag{A.29}$$

767 Substituting it back, yields

$$f(\widehat{\boldsymbol{y}}) \le f(\boldsymbol{x}) - \frac{1}{2L} \|\nabla f(\boldsymbol{x})\|^2, \quad \forall \boldsymbol{x} \in \mathbb{R}^n.$$
(A.30)

It follows that

$$f(\boldsymbol{x}^{\star}) \leq f(\widehat{\boldsymbol{y}}) \leq f(\boldsymbol{x}^{k}) - \frac{1}{2L} \|\nabla f(\boldsymbol{x}^{k})\|^{2}.$$
(A.31)
udes the proof.

Rearranging the terms concludes the proof.

A.7 PROOF OF PROPOSITION 3.3

Our Proposition 3.3, restated here as

Proposition A.5 (optimality gap). Let function $f : \mathbb{R}^n \to \mathbb{R}$ be L-smooth. Then,

$$\underbrace{\frac{\langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, -\nabla f(\boldsymbol{x}^{k}) \rangle}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}}_{optimal} - \underbrace{\frac{f(\boldsymbol{x}^{k}) - f(\boldsymbol{x}^{\star})}{\|\nabla f(\boldsymbol{x}^{k})\|^{2}}}_{estimated (\gamma_{0}=1)} \geq \frac{1}{2L}.$$
(A.32)

Proof. The proof follows instantly from (Nesterov, 2018, Theorem 2.1.5)

$$f(\boldsymbol{x}^{\star}) \geq f(\boldsymbol{x}^{k}) + \langle \boldsymbol{x}^{\star} - \boldsymbol{x}^{k}, \nabla f(\boldsymbol{x}^{k}) \rangle + \frac{1}{2L} \|\nabla f(\boldsymbol{x}^{\star}) - \nabla f(\boldsymbol{x}^{k})\|^{2}.$$
(A.33)

Dividing both sides with $\|
abla f(m{x}^k)\|^2$ (non-zero by assumption) concludes the proof.