

Theorem 2 proof

Anonymous Author(s)

Affiliation

Address

email

Abstract

1

2 **H1.** For any $t > 0$, the estimated parameter w_t stays within a ℓ_∞ -ball. There exists a constant
3 $W > 0$ such that $\|w_t\|_\infty \leq W$ almost surely.

4 **H2.** The function f is L -smooth (has L -Lipschitz gradients) w.r.t. the parameter w . There exists
5 some constant $L > 0$ such that for $(w, \vartheta) \in \Theta^2$, $f(w) - f(\vartheta) - \nabla f(\vartheta)^\top (w - \vartheta) \leq \frac{L}{2} \|w - \vartheta\|^2$.

6 We assume that the optimistic guess m_t at iteration t and the true gradient g_t are correlated:

7 **H3.** For any $t > 0$, $0 < \langle m_t | g_t \rangle = a_t \|g_t\|^2$ with some $0 < a_t \leq 1$, where $\langle | \rangle$ denotes the inner
8 product

9 We make a classical assumption in nonconvex optimization [?] on the magnitude of the gradient:

10 **H4.** There exists a constant $M > 0$ such that for any w and ξ , it holds $\|\nabla f(w, \xi)\| < M$.

11 **Lemma 1.** Assume H4, then the quantities defined in Algorithm ?? satisfy for any $w \in \Theta$ and $t > 0$,
12 $\|\nabla f(w_t)\| < M$, $\|\theta_t\| < M$ and $\|\hat{v}_t\| < M^2$.

13 **Lemma 2.** Assume H4, a strictly positive and a sequence of constant stepsizes $\{\eta_t\}_{t>0}$, $(\beta_1, \beta_2) \in$
14 $[0, 1]$, then the following holds:

$$\sum_{t=1}^{T_M} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} \theta_t \right\|_2^2 \right] \leq \frac{\eta^2 d T_M (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)}. \quad (1)$$

15 **Lemma 3.** Assume a strictly positive and non increasing sequence of stepsizes $\{\eta_t\}_{t>0}$, $\beta_1 < \beta_2 \in$
16 $[0, 1]$, then the following holds:

$$\bar{w}_{t+1} - \bar{w}_t \leq \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t,$$

17 where $\tilde{\theta}_t = \theta_t + \beta_1 \theta_{t-1}$ and $\tilde{g}_t = g_t - \beta_1 m_t + \beta_1 g_{t-1} + m_{t+1}$.

1 Proof of Theorem ??

19 **Proof** Using H2 and the iterate \bar{w}_t we have:

$$\begin{aligned} f(\bar{w}_{t+1}) &\leq f(\bar{w}_t) + \nabla f(\bar{w}_t)^\top (\bar{w}_{t+1} - \bar{w}_t) + \frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\|^2 \\ &\leq f(\bar{w}_t) + \underbrace{\nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t)}_A \\ &\quad + \underbrace{(\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t)}_B + \frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\|. \end{aligned} \quad (2)$$

20 **Term A.** Using Lemma 3, we have that:

$$\begin{aligned}\nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t) &\leq \nabla f(w_t)^\top \left[\frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \right] \\ &\leq \frac{\beta_1}{1 - \beta_1} \|\nabla f(w_t)\| \|\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2}\| \|\tilde{\theta}_{t-1}\| - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} \tilde{g}_t ,\end{aligned}$$

21 where the inequality is due to trivial inequality for positive diagonal matrix. Using Lemma 1 and
22 assumption H3 we obtain:

$$\nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t) \leq \frac{\beta_1(1 + \beta_1)}{1 - \beta_1} \mathbf{M}^2 [\|\eta_{t-1} \hat{v}_{t-1}^{-1/2}\| - \|\eta_t \hat{v}_t^{-1/2}\|] - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} \tilde{g}_t , \quad (3)$$

23 where we have used the fact that $\eta_t \hat{v}_t^{-1/2}$ is a diagonal matrix such that $\eta_{t-1} \hat{v}_{t-1}^{-1/2} \succcurlyeq \eta_t \hat{v}_t^{-1/2} \succcurlyeq 0$
24 (decreasing stepsize and max operator). Also note that:

$$\begin{aligned}-\nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} \tilde{g}_t &= -\nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} \bar{g}_t - \nabla f(w_t)^\top \left[\eta_t \hat{v}_t^{-1/2} - \eta_{t-1} \hat{v}_{t-1}^{-1/2} \right] \bar{g}_t \\ &\quad - \nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} (\beta_1 g_{t-1} + m_{t+1}) \\ &\leq -\nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} \bar{g}_t + (1 - a_t \beta_1) \mathbf{M}^2 [\|\eta_{t-1} \hat{v}_{t-1}^{-1/2}\| - \|\eta_t \hat{v}_t^{-1/2}\|] \\ &\quad - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} (\beta_1 g_{t-1} + m_{t+1}) ,\end{aligned} \quad (4)$$

25 where we have used Lemma 1 on $\|g_t\|$ and where that $\tilde{g}_t = \bar{g}_t + \beta_1 g_{t-1} + m_{t+1} = g_t - \beta_1 m_t +$
26 $\beta_1 g_{t-1} + m_{t+1}$. Plugging (4) into (3) yields:

$$\begin{aligned}\nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t) &\leq -\nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} \bar{g}_t + \frac{1}{1 - \beta_1} (a_t \beta_1^2 - 2a_t \beta_1 + \beta_1) \mathbf{M}^2 [\|\eta_{t-1} \hat{v}_{t-1}^{-1/2}\| - \|\eta_t \hat{v}_t^{-1/2}\|] \\ &\quad - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} (\beta_1 g_{t-1} + m_{t+1}) .\end{aligned} \quad (5)$$

27 **Term B.** By Cauchy-Schwarz (CS) inequality we have:

$$(\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t) \leq \|\nabla f(\bar{w}_t) - \nabla f(w_t)\| \|\bar{w}_{t+1} - \bar{w}_t\| . \quad (6)$$

28 Using smoothness assumption H2:

$$\begin{aligned}\|\nabla f(\bar{w}_t) - \nabla f(w_t)\| &\leq L \|\bar{w}_t - w_t\| \\ &\leq L \frac{\beta_1}{1 - \beta_1} \|w_t - \tilde{w}_{t-1}\| .\end{aligned} \quad (7)$$

29 By Lemma 3 we also have:

$$\begin{aligned}\bar{w}_{t+1} - \bar{w}_t &= \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \\ &= \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \eta_{t-1} \hat{v}_{t-1}^{-1/2} \left[I - (\eta_t \hat{v}_t^{-1/2})(\eta_{t-1} \hat{v}_{t-1}^{-1/2})^{-1} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \\ &= \frac{\beta_1}{1 - \beta_1} \left[I - (\eta_t \hat{v}_t^{-1/2})(\eta_{t-1} \hat{v}_{t-1}^{-1/2})^{-1} \right] (\tilde{w}_{t-1} - w_t) - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t ,\end{aligned} \quad (8)$$

30 where the last equality is due to $\tilde{\theta}_{t-1} \eta_{t-1} \hat{v}_{t-1}^{-1/2} = \tilde{w}_{t-1} - w_t$ by construction of $\tilde{\theta}_t$. Taking the
31 norms on both sides, observing $\|I - (\eta_t \hat{v}_t^{-1/2})(\eta_{t-1} \hat{v}_{t-1}^{-1/2})^{-1}\| \leq 1$ due to the decreasing stepsize
32 and the construction of \hat{v}_t and using CS inequality yield:

$$\|\bar{w}_{t+1} - \bar{w}_t\| \leq \frac{\beta_1}{1 - \beta_1} \|\tilde{w}_{t-1} - w_t\| + \|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\| . \quad (9)$$

We recall Young's inequality with a constant $\delta \in (0, 1)$ as follows:

$$\langle X | Y \rangle \leq \frac{1}{\delta} \|X\|^2 + \delta \|Y\|^2 .$$

33 Plugging (7) and (9) into (6) returns:

$$\begin{aligned} (\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t) &\leq L \frac{\beta_1}{1 - \beta_1} \|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\| \|w_t - \tilde{w}_{t-1}\| \\ &\quad + L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \|\tilde{w}_{t-1} - w_t\|^2. \end{aligned}$$

34 Applying Young's inequality with $\delta \rightarrow \frac{\beta_1}{1 - \beta_1}$ on the product $\|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\| \|w_t - \tilde{w}_{t-1}\|$ yields:

$$(\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t) \leq L \|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\|^2 + 2L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \|\tilde{w}_{t-1} - w_t\|^2. \quad (10)$$

35 The last term $\frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\|^2$ can be upper bounded using (9):

$$\begin{aligned} \frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\|^2 &\leq \frac{L}{2} \left[\frac{\beta_1}{1 - \beta_1} \|\tilde{w}_{t-1} - w_t\| + \|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\| \right] \\ &\leq L \|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\|^2 + 2L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \|\tilde{w}_{t-1} - w_t\|^2. \end{aligned} \quad (11)$$

36 Plugging (5), (10) and (11) into (2) and taking the expectations on both sides give:

$$\begin{aligned} &\mathbb{E} \left[f(\bar{w}_{t+1}) + \frac{1}{1 - \beta_1} \tilde{M}_t^2 \|\eta_t \hat{v}_t^{-1/2}\| - \left(f(\bar{w}_t) + \frac{1}{1 - \beta_1} \tilde{M}_t^2 \|\eta_{t-1} \hat{v}_{t-1}^{-1/2}\| \right) \right] \\ &\leq \mathbb{E} \left[-\nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} \tilde{g}_t - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} (\beta_1 g_{t-1} + m_{t+1}) \right] \\ &\quad + \mathbb{E} \left[2L \|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\|^2 + 4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \|\tilde{w}_{t-1} - w_t\|^2 \right], \end{aligned}$$

37 where $\tilde{M}_t^2 = (a_t \beta_1^2 + \beta_1) M^2$. Note that the expectation of \tilde{g}_t conditioned on the filtration \mathcal{F}_t reads
38 as follows

$$\mathbb{E} [\nabla f(w_t)^\top \tilde{g}_t] = \mathbb{E} [\nabla f(w_t)^\top (g_t - \beta_1 m_t)] = (1 - a_t \beta_1) \|\nabla f(w_t)\|^2. \quad (12)$$

39 Summing from $t = 1$ to $t = T$ leads to

$$\begin{aligned} &\frac{1}{M} \sum_{t=1}^{T_M} ((1 - a_t \beta_1) \eta_{t-1} + (\beta_1 + a_t) \eta_t) \|\nabla f(w_t)\|^2 \leq \\ &\mathbb{E} \left[f(\bar{w}_1) + \frac{1}{1 - \beta_1} \tilde{M}_t^2 \|\eta_0 \hat{v}_0^{-1/2}\| - \left(f(\bar{w}_{T_M+1}) + \frac{1}{1 - \beta_1} \tilde{M}_t^2 \|\eta_{T_M} \hat{v}_{T_M}^{-1/2}\| \right) \right] \\ &\quad + 2L \sum_{t=1}^{T_M} \mathbb{E} [\|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\|^2] + 4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \sum_{t=1}^{T_M} \mathbb{E} [\|\tilde{w}_{t-1} - w_t\|^2] \\ &\leq \mathbb{E} \left[\Delta f + \frac{1}{1 - \beta_1} \tilde{M}_t^2 \|\eta_0 \hat{v}_0^{-1/2}\| \right] + 2L \sum_{t=1}^{T_M} \mathbb{E} [\|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\|^2] \\ &\quad + 4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \sum_{t=1}^{T_M} \mathbb{E} [\|\tilde{w}_{t-1} - w_t\|^2], \end{aligned} \quad (13)$$

40 where we denote $\Delta f := f(\bar{w}_1) - f(\bar{w}_{T_M+1})$. We note that by definition of \hat{v}_t , and a constant
41 learning rate η_t , we have

$$\begin{aligned} \|\tilde{w}_{t-1} - w_t\|^2 &= \|\eta_{t-1} \hat{v}_{t-1}^{-1/2} (\theta_{t-1} + h_t)\|^2 \\ &= \|\eta_{t-1} \hat{v}_{t-1}^{-1/2} (\theta_{t-1} + \beta_1 \theta_{t-2} + (1 - \beta_1) m_t)\|^2 \\ &\leq \|\eta_{t-1} \hat{v}_{t-1}^{-1/2} \theta_{t-1}\|^2 + \|\eta_{t-2} \hat{v}_{t-2}^{-1/2} \beta_1 \theta_{t-2}\|^2 + (1 - \beta_1)^2 \|\eta_{t-1} \hat{v}_{t-1}^{-1/2} m_t\|^2. \end{aligned}$$

42 Using Lemma 2 we have

$$\begin{aligned} & \sum_{t=1}^{T_M} \mathbb{E} [\|\tilde{w}_{t-1} - w_t\|^2] \\ & \leq (1 + \beta_1^2) \frac{\eta^2 d T_M (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} + (1 - \beta_1)^2 \sum_{t=1}^{T_M} \mathbb{E} [\|\eta_{t-1} \hat{v}_{t-1}^{-1/2} m_t\|] . \end{aligned}$$

43 Here, if we assume $\min_{1, \dots, T_M} a_t = a_m$ and constant $\text{lr } \eta$, we have from (13)

$$\frac{\eta}{M} \sum_1^{T_M} ((1 + \beta_1) + (1 - \beta_1) a_t) \|\nabla f(w_t)\|^2 \geq \frac{\eta}{M} ((1 + \beta_1) + (1 - \beta_1) a_m) \sum_1^{T_M} \|\nabla f(w_t)\|^2$$

44 Then

$$\mathbb{E} [\|\nabla f(w_T)\|^2] = \frac{\sum_1^{T_M} \|\nabla f(w_t)\|^2}{T_M} \leq \frac{M}{T_M \eta ((1 + \beta_1) + (1 - \beta_1) a_m)} \times (13) RHS \leq \dots$$

45 Say if η is still $\frac{1}{\sqrt{dT_M}}$, the order of the bound should be the same.

46 Assume $a_m = \min_{1, \dots, T_M} a_t$ and denote $\tilde{M}_m^2 = (a_m \beta_1^2 + \beta_1) M^2$. Setting a constant learning rate
47 $\eta_t = \eta$ and plugging in (13) yields:

$$\begin{aligned} \mathbb{E} [\|\nabla f(w_T)\|^2] &= \frac{1}{\sum_{j=1}^{T_M} \eta_j} \sum_{t=1}^{T_M} \eta_t \|\nabla f(w_t)\|^2 = \frac{\sum_1^{T_M} \|\nabla f(w_t)\|^2}{T_M} \\ &\leq \frac{M}{T_M \eta ((1 - a_m \beta_1) + (\beta_1 + a_m))} \mathbb{E} \left[\Delta f + \frac{1}{1 - \beta_1} \tilde{M}_m^2 \|\eta_0 \hat{v}_0^{-1/2}\| \right] \\ &\quad + \frac{4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 M}{T_M \eta ((1 - a_m \beta_1) + (\beta_1 + a_m))} (1 + \beta_1^2) \frac{\eta^2 d T_M (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} \\ &\quad + \frac{M}{T_M \eta ((1 - a_m \beta_1) + (\beta_1 + a_m))} (1 - \beta_1)^2 \sum_{t=1}^{T_M} \mathbb{E} [\|\eta_{t-1} \hat{v}_{t-1}^{-1/2} m_t\|] \\ &\quad + \frac{2LM}{T_M \eta ((1 - a_m \beta_1) + (\beta_1 + a_m))} \sum_{t=1}^{T_M} \mathbb{E} [\|\eta_t \hat{v}_t^{-1/2} \tilde{g}_t\|^2] , \end{aligned}$$

48 where T_M is a random termination number distributed according (??).

49 Setting the stepsize to $\eta = \frac{1}{\sqrt{dT_M}}$ yields :

$$\mathbb{E} [\|\nabla f(w_T)\|^2] \leq C_{1,m} \sqrt{\frac{d}{T_M}} + C_{2,m} \frac{1}{T_M} + \frac{\eta}{T_M} D_{1,m} \mathbb{E} [\|\hat{v}_{t-1}^{-1/2} m_t\|] + \frac{\eta}{T_M} D_{2,m} \mathbb{E} [\|\hat{v}_{t-1}^{-1/2} \tilde{g}_t\|] ,$$

50 where

$$\begin{aligned} C_{1,m} &= \frac{M}{(1 - a_m \beta_1) + (\beta_1 + a_m)} \Delta f + \frac{4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 M}{(1 - a_m \beta_1) + (\beta_1 + a_m)} \frac{(1 + \beta_1^2)(1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} , \\ C_{2,m} &= \frac{M}{(1 - \beta_1) ((1 - a_m \beta_1) + (\beta_1 + a_m))} (a_m \beta_1^2 + \beta_1) M^2 \mathbb{E} [\|\hat{v}_0^{-1/2}\|] . \end{aligned}$$

51 **Simple case as in [?]:** if $\beta_1 = 0$ then $\tilde{g}_t = g_t + m_{t+1}$ and $g_t = \theta_t$. Also using Lemma 2 we have
52 that:

$$\sum_{t=1}^{T_M} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} g_t \right\|_2^2 \right] \leq \frac{\eta^2 d T_M}{(1 - \beta_2)} ;$$

53 which leads to the final bound:

$$\mathbb{E} [\|\nabla f(w_T)\|^2] \leq \sqrt{\frac{d}{T_M}} \tilde{C}_{1,m} + \frac{1}{T_M} \tilde{C}_{2,m} ,$$

54 where

$$\begin{aligned}\tilde{C}_{1,m} &= C_{1,m} + \frac{\mathbf{M}}{(1 - a_m\beta_1) + (\beta_1 + a_m)} \left[\frac{a_m(1 - \beta_1)^2}{1 - \beta_2} + 2L \frac{1}{1 - \beta_2} \right] , \\ \tilde{C}_{2,m} &= C_{2,m} = \frac{\mathbf{M}}{(1 - \beta_1) ((1 - a_m\beta_1) + (\beta_1 + a_m))} \tilde{\mathbf{M}}_m^2 \mathbb{E}[\|\hat{v}_0^{-1/2}\|] .\end{aligned}$$

55

□