

ADAPTIVE LEARNING RATES WITH MAXIMUM VARIATION AVERAGING

Anonymous authors

Paper under double-blind review

ABSTRACT

Adaptive gradient methods such as RMSPROP and ADAM use exponential moving estimate of the squared gradient to compute coordinate-wise adaptive step sizes, achieving better convergence than SGD in face of noisy objectives. However, ADAM can have undesirable convergence behaviors due to unstable or extreme adaptive learning rates. Methods such as AMSGRAD and ADABOUND have been proposed to stabilize the adaptive learning rates of ADAM in the later stage of training, but they do not outperform ADAM in some practical tasks such as training Transformers (Vaswani et al., 2017). In this paper, we propose an adaptive learning rate principle, in which the running mean of squared gradient is replaced by a weighted mean, with weights chosen to maximize the estimated variance of each coordinate. This gives a worst-case estimate for the local gradient variance, taking smaller steps when large curvatures or noisy gradients are present, which leads to more desirable convergence behaviors than ADAM. We prove the proposed algorithm converges under mild assumptions for nonconvex stochastic optimization problems, and demonstrate the improved efficacy of our adaptive averaging approach on image classification, machine translation and natural language understanding tasks. Moreover, our method overcomes the non-convergence issue of ADAM in BERT pretraining at large batch sizes, while achieving better test performance than LAMB in the same setting.

1 INTRODUCTION

Stochastic Gradient Descent (SGD) and its variants are commonly used for training deep neural networks because of their effectiveness and efficiency. In their simplest form, gradient methods train a network by iteratively moving each parameter in the direction of the negative gradient (or the running average of gradients) of the loss function on a randomly sampled mini-batch of training data. A scalar learning rate is also applied to control the size of the update. In contrast, *adaptive* stochastic gradient methods use coordinate-specific learning rates, which are inversely proportional to the square root of the running mean of squared gradients (Tieleman & Hinton, 2012; Duchi et al., 2011; Kingma & Ba, 2015). Such methods are proposed to improve the stability of SGD on non-stationary problems, and have achieved success in different fields across Speech, Computer Vision (CV), and Natural Language Processing (NLP).

Despite the popularity of adaptive methods such as ADAM (Kingma & Ba, 2015), the instability of their adaptive learning rates sometimes leads to sub-optimal solutions or even non-convergent behavior on some simple problems (Reddi et al., 2018; Luo et al., 2019). AMSGRAD (Reddi et al., 2018) was proposed to stabilize ADAM by computing the adaptive learning rate with an update rule that guarantees monotonically decaying adaptive learning rates for each coordinate. ADABOUND (Luo et al., 2019) clips the adaptive learning rate of ADAM with a decreasing upper bound and an increasing lower bound, so that it transitions into SGD in the final stage of training. However, to our knowledge, neither of the two approaches have been deployed to enhance ADAM on recent large-scale problems such as training Transformer-based language models (Devlin et al., 2019; Liu et al., 2019; Lan et al., 2020; Raffel et al., 2019; Zhu et al., 2020). The stochastic gradients of Transformer’s loss functions exhibit heavy-tailed statistics, making SGD converge to a much worse solution than ADAM. Zhang et al. (2019b) has to use gradient clipping to stabilize SGD for training transformers, indicating that the strategy of ADABOUND, which is to transition from ADAM into SGD, will fail on Transformers (see Appendix D for instance). RAdam (Liu et al., 2020) was re-

cently invented to free ADAM from the warmup schedule for training Transformers, but its variance rectification term does not really depend on the observed gradients during training, and Ma & Yarats (2019) found that using a linear warmup over $2 \cdot (1 - \beta_2)^{-1}$ iterations for ADAM achieves almost the same convergence as RAdam.

In this work, we explore a different approach to improving the stability of adaptive learning rates. We propose *Maximum Variation Averaging* (MaxVA), which computes the running average of squared gradients using dynamic, rather than constant, coordinate-wise weights. These weights are chosen so that the estimated variance of gradients is maximized. The MaxVA weights for maximizing this variance have a simple closed-form solution that requires little storage or computational cost. With this solution, MaxVA assigns a higher weight to a coordinate if the gradient at that coordinate deviates too much (compared with the estimated variance) from the estimated mean, enabling a faster adaptation to gradient change or the curvature. Extensive experiments on synthetic and practical datasets demonstrate that this leads to an improved adaptability and stability for ADAM, yielding better test set performance than ADAM on a variety of tasks. The effect is especially evident in the large-batch optimization setting for BERT, where the total number of iterations is sharply reduced and a faster adaptation in each step is more important, and ADAM with MaxVA converges faster than ADAM and LAMB (You et al., 2020) with higher test scores in the experiments. We also prove MaxVA converges under mild assumptions in the nonconvex stochastic optimization setting.

2 PRELIMINARY AND DEFINITIONS

By default, all vector-vector operators are element-wise in the following sections. Let $\theta \in \mathbb{R}^d$ be the parameters of the network to be trained, $\ell(x; \theta)$ is the loss of the model with parameters θ evaluated at x . Our goal is to minimize the expected risk on the data distribution defined as:

$$f(\theta) = \mathbb{E}_{x \sim \mathcal{D}} [\ell(x; \theta)]. \quad (1)$$

In most deep learning problems, only a finite number of potentially noisy samples can be used to approximate Eq. 1, and the gradients are computed on randomly sampled minibatches during training. Stochastic regularizations such as Dropout (Srivastava et al., 2014) are commonly used for training Transformers (Vaswani et al., 2017), which further adds to the randomness of the gradients. Thus, it is important to design optimizers that tolerate noisy gradients. ADAM (Kingma & Ba, 2015) is an effective optimizer that adapts to such noisy gradients. It keeps exponential moving averages to estimate the first and second moment of the gradient, m_t and v_t , defined as:

$$\begin{aligned} \tilde{m}_t &= \alpha \tilde{m}_{t-1} + (1 - \alpha) g_t, & m_t &= \frac{\tilde{m}_t}{1 - \alpha^{t+1}}, \\ \tilde{v}_t &= \beta \tilde{v}_{t-1} + (1 - \beta) g_t^2, & v_t &= \frac{\tilde{v}_t}{1 - \beta^{t+1}}, \end{aligned} \quad (2)$$

where $\alpha, \beta \in [0, 1]$, $g_t = \nabla_{\theta} \ell(x_t; \theta_t)$ is the gradient of the t -th minibatch x_t , $\tilde{m}_0 = \tilde{v}_0 = 0$, and m_t, v_t corrects the initialization bias of \tilde{m}_t, \tilde{v}_t (Kingma & Ba, 2015). ADAM updates the parameters with the estimated moments as $\theta_{t+1} = \theta_t - \eta_t \frac{m_t}{\sqrt{v_t + \epsilon}}$, where $\epsilon > 0$ is a small constant for numerical stability.

If we assume that the distribution of the stochastic gradient is constant within the effective horizon of the running average (Balles & Hennig, 2018), an assumption that is more accurate when the model is closer to convergence, then m_t and v_t will be unbiased estimates of the first and second moments of the gradient g_t . Same as other adaptive methods such as ADAM and the recently proposed AdaBelief (Zhuang et al., 2020), we use this assumption throughout training. Specifically, let σ_t^2 be the variance of g_t . At time t , we assume $\mathbb{E}[m_t] \approx \nabla f_t$, $\mathbb{E}[v_t] \approx \nabla f_t^2 + \sigma_t^2$. ADAM, RMSPROP and other variants that divide the update steps by $\sqrt{v_t}$ can be seen as adapting to the gradient variance under this assumption. These adaptive methods take smaller step sizes when the estimated variance $\tilde{\sigma}_t^2 = v_t - m_t^2$ is high. Higher local gradient variance indicates higher local curvature, and vice versa. In certain quadratic approximations to the loss function, this variance is proportional to the curvature (Schaul et al., 2013) (Eq. 13 of our paper). Therefore, like a diagonal approximation to Newton’s method, using such adaptative learning rates to adapt to the curvature can accelerate the convergence of first-order methods.

However, the adaptive learning rate $\eta_t/(\sqrt{v_t} + \epsilon)$ of ADAM and RMSPROP can take extreme values, causing convergence to undesirable solutions (Wilson et al., 2017; Chen et al., 2019). Reddi et al. (2018) gave one such counter example where gradients in the correct direction are large but occur at a low frequency, and ADAM converges to the solution of maximum regret. They solve this issue by keeping track of the maximum v_t for each coordinate throughout training with a new variable \hat{v}_t , and replace the adaptive learning rate with $\eta_t/\sqrt{\hat{v}_t}$ to enforce monotonically decreasing learning rates. Extremely small adaptive learning rates can also cause undesirable convergence behavior, as demonstrated by a counter example from Luo et al. (2019).

3 MAXIMIZING THE VARIANCE OF RUNNING ESTIMATIONS

Algorithm 1 MADAM

1: **Input:** Learning rate $\{\eta_t\}_{t=1}^T$, parameter $0 < \alpha < 1, 0 < \beta < \bar{\beta} < 1, \epsilon > 0$
2: Set $\tilde{m}_0 = \tilde{u}_0 = \tilde{v}_0 = w_0 = 0$
3: **for** $t = 1$ **to** T **do**
4: Draw samples S_t from training set
5: Compute $g_t = \frac{1}{|S_t|} \sum_{x_k \in S_t} \nabla \ell(x_k; \theta_t)$
6: $\tilde{m}_t = \alpha \tilde{m}_{t-1} + (1 - \alpha) g_t$
7: $\tilde{\beta}_t = \arg \max_{\beta} v_t(\beta) - u_t^2(\beta) \quad \triangleright$ see Eq 7
8: $\beta_t = \max(\beta, \min(\bar{\beta}, \tilde{\beta}_t))$
9: $\tilde{u}_t = \beta_t \tilde{u}_{t-1} + (1 - \beta_t) g_t$
10: $\tilde{v}_t = \beta_t \tilde{v}_{t-1} + (1 - \beta_t) g_t^2$
11: $w_t = \beta_t w_{t-1} + (1 - \beta_t)$
12: $\theta_t = \theta_{t-1} - \eta_t \frac{\sqrt{w_t}}{1 - \alpha^t} \frac{\tilde{m}_t}{\sqrt{\tilde{v}_t} + \epsilon}$

Algorithm 2 LAMADAM

1: **Input:** Learning rate $\{\eta_t\}_{t=1}^T$, parameter $0 < \alpha < 1, 0 < \beta < \bar{\beta} < 1, \epsilon > 0$
2: Set $\tilde{m}_0 = \tilde{u}_0 = \tilde{v}_0 = w_0 = 0$
3: **for** $t = 1$ **to** T **do**
4: Draw samples S_t from training set
5: Compute $g_t = \frac{1}{|S_t|} \sum_{x_k \in S_t} \nabla \ell(x_k; \theta_t)$
6: $\tilde{\beta}_t = \arg \max_{\beta} v_t(\beta) - u_t^2(\beta) \quad \triangleright$ see Eq 7
7: $\beta_t = \max(\beta, \min(\bar{\beta}, \tilde{\beta}_t))$
8: $\tilde{u}_t = \beta_t \tilde{u}_{t-1} + (1 - \beta_t) g_t$
9: $\tilde{v}_t = \beta_t \tilde{v}_{t-1} + (1 - \beta_t) g_t^2$
10: $w_t = \beta_t w_{t-1} + (1 - \beta_t)$
11: $\tilde{m}_t = \alpha \tilde{m}_{t-1} + (1 - \alpha) \frac{g_t}{\sqrt{\tilde{v}_t/w_t} + \epsilon}$
12: $\theta_t = \theta_{t-1} - \frac{\eta_t}{1 - \alpha^t} \tilde{m}_t$

Motivation. We propose to mitigate the undesirable convergence issue of ADAM by changing the constant running average coefficient β for the second moment into an adaptive one. The idea is to allow β_t to adopt the value that maximizes the estimated variance of the gradient at each iteration t . Therefore, our algorithm can use β_t as the adaptive running average coefficient to take steps that are conservative enough to avoid instability but aggressive enough to make progress.

Maximum Variation Averaging. Formally, we estimate the variance of the gradient at each coordinate by keeping track of the zeroth, first, and second moments of the gradient as functions of the adaptive running average coefficient β_t , denoted as $w_t(\beta_t)$, $\tilde{u}_t(\beta_t)$ and $\tilde{v}_t(\beta_t)$, respectively:

$$w_t(\beta_t) = \beta_t w_{t-1}(\beta_{t-1}) + (1 - \beta_t), \quad (3)$$

$$\tilde{u}_t(\beta_t) = \beta_t \tilde{u}_{t-1}(\beta_{t-1}) + (1 - \beta_t) g_t, \quad (4)$$

$$\tilde{v}_t(\beta_t) = \beta_t \tilde{v}_{t-1}(\beta_{t-1}) + (1 - \beta_t) g_t^2. \quad (5)$$

The zeroth moment $w_t(\beta_t)$ is used to normalize $\tilde{u}_t(\beta_t)$ and $\tilde{v}_t(\beta_t)$ to achieve bias-corrected estimates $u_t(\beta_t) = \tilde{u}_t(\beta_t)/w_t(\beta_t)$ and $v_t(\beta_t) = \tilde{v}_t(\beta_t)/w_t(\beta_t)$ for the first and second moments, so that the estimates are not biased towards $\tilde{m}_0 = \tilde{v}_0 = 0$ (Kingma & Ba, 2015).

Under our assumptions, the bias-corrected local estimate of the gradient variance is $\tilde{\sigma}_t^2 = \tilde{v}_t(\beta_t)/w_t(\beta_t) - [\tilde{u}_t(\beta_t)/w_t(\beta_t)]^2$. Taking the $\arg \max$ for $\tilde{\sigma}_t^2$, we find the β_t that achieves the worst-case (maximal) variance for each coordinate i :

$$\beta_{t,i} = \arg \max_{\beta} \tilde{\sigma}_{t,i}^2 = \arg \max_{\beta} v_{t,i}(\beta) - [u_{t,i}(\beta)]^2. \quad (6)$$

We call our approach to finding adaptive running average coefficient β_t *Maximum Variation Averaging* (MaxVA). We plug MaxVA into ADAM and its variant LAPROP (Ziyin et al., 2020), which results in two novel algorithms, MADAM and LAMADAM, listed in Algorithm 1 and Algorithm 2. Different from ADAM, LAPROP uses v_t to normalize the gradients before taking the running average, which results in higher empirical stability under various hyperparameters. Note, we only use the MaxVA formula for the *second* moment $u_t(\beta_t)$ used for scaling the learning rate; m_t is still an exponential moving average (with a constant coefficient α) of the gradient for MADAM or the normalized gradient for LAMADAM.

Finding β_t via a Closed-form Solution. The maximization for β_t in Eq. 6 is quadratic and has a relatively simple closed-form solution that produces maximal $\tilde{\sigma}_t^2$ for each coordinate:

$$\beta_t = \frac{(g_t - u_{t-1})^2 + v_{t-1} - u_{t-1}^2}{w_{t-1}((g_t - u_{t-1})^2 - v_{t-1} + u_{t-1}^2) + (g_t - u_{t-1})^2 + v_{t-1} - u_{t-1}^2}, \quad (7)$$

where all variables are vectors and all the operations are elementwise, and we have abbreviated $u_{t-1}(\beta_{t-1})$, $v_{t-1}(\beta_{t-1})$ and $w_{t-1}(\beta_{t-1})$ into u_{t-1} , v_{t-1} and w_{t-1} . We use this abbreviation in the following sections, and defer the derivation of Eq. 7 to Appendix A.

Implementation Notes. We apply MaxVA in every step except for the first step, where the gradient variance one can observe is zero. So for Algorithm 1 and Algorithm 2 we define:

$$\tilde{u}_1 = (1 - \beta_1)g_1, \tilde{v}_1 = (1 - \beta_1)g_1^2, w_1 = 1 - \beta_1. \quad (8)$$

The coefficient β_1 for $t = 1$ is set to a constant that is the same as typical values for ADAM. To obtain a valid running average, we clip β_t so that $\underline{\beta} \leq \beta_t \leq \bar{\beta}$, where the typical values are $\underline{\beta} = 0.5, 0.98 \leq \bar{\beta} \leq 1$. For convenience, we set $\beta_1 = \bar{\beta}$ by default. For $t > 1$, since $0 < \beta_t \leq 1$, w_t will monotonically increase from $(1 - \beta_1)$ to 1. Before clipping, for any g_t, u_{t-1}, v_{t-1} satisfying $v_{t-1} - u_{t-1}^2 > 0$ in Eq. 7, we have $\beta_t \in [1/(1 + w_{t-1}), 1/(1 - w_{t-1})]$. As a result, the lower bound that we use ($\underline{\beta} = 0.5$) is tight and does not really change the value of β_t , and as $t \rightarrow \infty$, $w_t \rightarrow 1$ and $\beta_t \in [0.5, \infty]$. We have a special case at $t = 2$, where β_t is a constant $1/(2 - \beta_1)$.

In practice, we also add a small coefficient $\delta > 0$ to the denominator of Eq. 7 to prevent division by zero, which will have negligible effect on the value of β_t and does not violate the maximum variation objective (Eq. 6). All the derivations for these conclusions are deferred to Appendix C.

Effect of Maximum Variation Averaging. By definition, we have $v_{t-1} - u_{t-1}^2 \geq 0$, but in most cases $v_{t-1} - u_{t-1}^2 > 0$. If we define a new variable $R_t = (g_t - u_{t-1})^2 / (v_{t-1} - u_{t-1}^2)$, which represents the degree of deviation of gradient g_t from the current estimated average, we can rewrite:

$$\beta_t = \frac{R_t + 1}{(1 + w_t)R_t + 1 - w_t}. \quad (9)$$

From Eq. 9, we can see β_t monotonically decreases from $1/(1 - w_t)$ to $1/(1 + w_t)$ as R_t increases from 0 to ∞ , and equals to 1 when $R_t = 1$. As a result, for each entry, if $R_t \leq 1$, or the deviation of the gradient g_t from the current running mean u_{t-1} is within the estimated standard deviation $\tilde{\sigma}$, we will use β to update \tilde{v}_t , which is the smallest change we allow for \tilde{v}_t . If g_t deviates much more than $\tilde{\sigma}$ from u_t , MaxVA will find a smaller β_t and therefore a higher weight $(1 - \beta_t)$ on g_t^2 to adapt to the change faster. For example, when an abnormally large gradient entry occurs in some coordinate i , MaxVA will assign a smaller $\beta_{t,i}$ to obtain a larger $v_{t,i}$ and smaller adaptive step size compared with ADAM with $\beta = \bar{\beta}$. This allows a quick response to impede abnormally large gradients, which enables a better handling for the heavy-tailed distribution of gradients in the process of training Transformers (Zhang et al., 2019b). As a side effect, \tilde{v}_t tends to be larger than ADAM/LAPROP using a constant β , but as we will show in the experiments, using a larger learning rate counters such an effect and achieves better results.

On the other hand, when the variance decreases in the later phase of training, g_t tends to be within $\tilde{\sigma}_t$, and MaxVA tends to find the slowest rate for decreasing \tilde{v}_t . This allows large values of \tilde{v}_t to last for a longer horizon even compared with setting β_t to a constant $\bar{\beta}$ on the same sequence, since we have assigned more mass to large gradients, which can be seen as an adaptive version of AMSGRAD. Note that MaxVA and AMSGRAD can be complementary approaches if applied together, which we have found helpful for Image Classification on CIFAR10/100.

Convergence Analysis. We prove the convergence of MaxVA in the nonconvex stochastic optimization setting. For the sake of simplicity, we analyze the case where $\alpha = 0$, which is effectively applying MaxVA to RMSPROP. We leave the analysis for $\alpha \neq 0$ for future research. We assume the function ℓ is L -smooth in θ , i.e., there exists a constant L such that

$$\|\nabla \ell(x; \theta_1) - \nabla \ell(x; \theta_2)\| \leq L \|\theta_1 - \theta_2\|, \text{ for all } \theta_1, \theta_2 \in \mathbb{R}^d, x \in \mathcal{X}. \quad (10)$$

This automatically implies that $f(\theta) = \mathbb{E}[\ell(x; \theta)]$ is L -smooth. Such a smoothness assumption holds for networks with smooth activation functions, e.g., Transformers that use the GELU activation (Hendrycks & Gimpel, 2016). We also need to assume function ℓ has bounded gradient, i.e.,

$\|\nabla_{\theta}\ell(x; \theta)\|_{\infty} \leq G$ for all $\theta \in \mathbb{R}^d, x \in \mathcal{X}$. As typically used in the analysis of stochastic first-order methods (Zaheer et al., 2018; Ghadimi & Lan, 2013), we assume the stochastic gradient has bounded variance: $\mathbb{E}[\|\nabla_{\theta}\ell(x; \theta)_i - \nabla_{\theta}f(\theta)_i\|^2] \leq \sigma^2$ for all $\theta \in \mathbb{R}^d$. Further, we assume the batch size increases with time as $b_t = t$, which is also adopted in the analysis of SIGNSGD (Bernstein et al., 2018), and holds in our large batch experiments. Theorem 1 gives a “worst-case” convergence rate of MaxVA to a stationary point under these assumptions, where the dependence of β_t on g_t is ignored and we only consider the worst-case of β_t in each step. MADAM still converges under this setting. The proof is given in Appendix B.

Theorem 1. Define $w_0 = 1$. Let $\eta_t = \eta$ and $b_t = t$ for all $t \in [T]$. Furthermore, we assume $\epsilon, \beta, \bar{\beta}, \eta$ are chosen such that $\eta \leq \frac{\epsilon}{2L}$, $1 - \underline{\beta} \leq \frac{\epsilon^2}{16G^2}$, and $\bar{\beta} \leq 2\underline{\beta}$. Then for θ_t generated using MADAM, we have the following bound:

$$\mathbb{E}\|\nabla f(\theta_a)\|^2 \leq O\left(\frac{f(\theta_1) - f(\theta^*)}{\eta T} + \frac{2\sigma dG}{\epsilon\sqrt{T}}\right), \quad (11)$$

where θ^* is an optimal solution to minimize the objective in Eq. 1, and θ_a is an iterate uniformly randomly chosen from $\{\theta_1, \dots, \theta_T\}$.

4 EXPERIMENTS ON SYNTHETIC DATA

For a better control of data distribution and demonstrate the efficacy of MaxVA with statistical significance on a larger number of instances, we compare MADAM and the baselines in two sets of synthetic data. The first dataset simulates prevalent machine learning settings where mini-batch stochastic gradient methods are applied on a finite set of samples, on which we show MADAM fixes the nonconvergence issue of ADAM and achieves faster convergence rate than AMSGRAD. The second dataset evaluates the algorithms under different curvatures and gradient noise levels, where we show MADAM achieves both lower loss and variance than fine-tuned ADAM at convergence.

4.1 CONVERGENCE WITH STOCHASTIC GRADIENTS

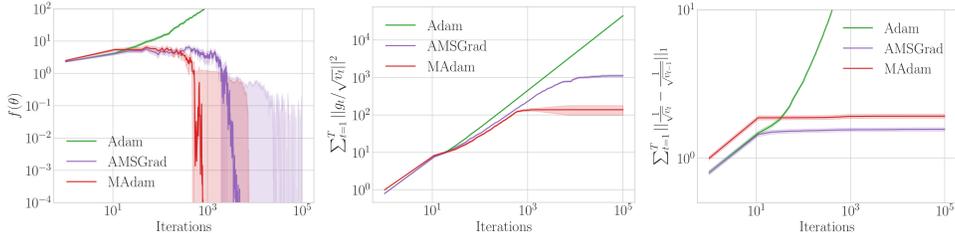


Figure 1: Median and standard error (over 100 runs) of objective value ($f(\theta)$), accumulated update size ($\sum_{t=1}^T \|g_t/\sqrt{v_t}\|^2$) and total change in adaptive learning rate ($\sum_{t=1}^T \|\frac{1}{\sqrt{v_t}} - \frac{1}{\sqrt{v_{t-1}}}\|_1$) for ADAM, AMSGRAD, MADAM on the problem defined in Eq. 12.

Since MaxVA maximizes the variance and the gradient converges to zero in most cases, MADAM biases towards larger v_t than ADAM but does not require v_t to be monotonically increasing, which is like an adaptive version of AMSGRAD. To highlight the difference, we compare ADAM, MADAM and AMSGRAD on the synthetic dataset from Chen et al. (2019) simulating training with stochastic mini batches on a finite set of samples. Formally, let $\mathbb{I}[\cdot]$ be the indicator function. We consider the problem $\min_{\theta} f(\theta) = \sum_{i=1}^{11} \ell_i(\theta)$ where

$$\ell_i(\theta) = \begin{cases} \mathbb{I}[i=1]5.5\theta^2 + \mathbb{I}[i \neq 1](-0.5\theta^2), & \text{if } |\theta| \leq 1; \\ \mathbb{I}[i=1](11|\theta| - 5.5) + \mathbb{I}[i \neq 1](-|\theta| + 0.5), & \text{otherwise.} \end{cases} \quad (12)$$

At every step, a random index i is sampled uniformly from $i \in [11]$, and the gradient $\nabla \ell_i(\theta)$ is used by the optimizer. The only stationary point where $\nabla f(\theta) = 0$ is $\theta = 0$. We set $\alpha = 0, \beta = 0.9$ for ADAM and AMSGRAD. For MADAM, we set $\alpha = 0, (\beta, \bar{\beta}) = (0.5, 1)$. We select the best constant learning rates for the three algorithms, see Appendix E for details.

We plot the median and standard error of the objective ($f(\theta)$), accumulated update size ($S_1 = \sum_{t=1}^T \|g_t/\sqrt{v_t}\|^2$), and total change in adaptive step size ($S_2 = \sum_{t=1}^T \|\frac{1}{\sqrt{v_t}} - \frac{1}{\sqrt{v_{t-1}}}\|_1$) over 100 runs in Figure 12. The optimal learning rates for these optimizers are different, so for fair comparisons, we have ignored the constant learning rate in S_1 and S_2 . From the curves of $f(\theta)$, we can see ADAM diverges, and MADAM converges faster than AMSGRAD in the later stage. As shown by the S_2 curves, the adaptive step sizes of MADAM and AMSGRAD all converged to some constant values after about 10 steps, but MADAM converges faster on both $f(\theta)$ and S_1 , indicating the adaptive step size found by MADAM fits the geometry of the problem better than AMSGRAD. This also shows $S_1 + S_2$ of MADAM has a smaller slope than AMSGRAD in the log-scale plots after 10 iterations, leading to a faster theoretical convergence rate in the bound given by Chen et al. (2019). The slightly larger variation in adaptive step sizes of MADAM at the beginning of training, shown by the larger S_2 values, demonstrates MADAM adapts faster to the changing gradients than AMSGRAD, achieved by dynamically selecting $\beta < 0.9$.

4.2 CONVERGENCE IN THE NOISY QUADRATIC MODEL

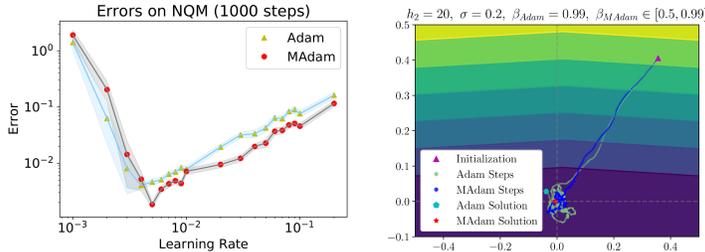


Figure 2: Results on NQM. The left figure shows the mean and standard error of the loss under different learning rates η , computed over 100 runs at each point. We select the best β for ADAM at each η . The best results (mean and variance) of ADAM and MADAM are $1.84e-3$ ($2.51e-4$) and $4.05e-3$ ($4.84e-4$) respectively. Figure on the right gives a qualitative example of the trajectories of two approaches.

We analyze the ability of MADAM to adapt to curvature and gradient noise on the simple but illustrative Noisy Quadratic Model (NQM), which has been widely adopted for analyzing optimization dynamics (Schaul et al., 2013; Wu et al., 2018; Zhang et al., 2019a;c). The loss function is defined as $f(\theta) = \mathbb{E}_{x \sim \mathcal{N}(0, \sigma^2 I)} \left[\frac{1}{2} \sum_{i=1}^d h_i (\theta_i - x_i)^2 \right]$, where x is a noisy observation of the ground-truth parameter $\theta^* = 0$, simulating the gradient noise in stochastic optimization, and h_i represents the curvature of the system in d dimensions. In each step, the following noisy gradient for each coordinate i are put into the algorithms, from which we can see the gradient noise is proportional to the curvature h_i :

$$\nabla_{\theta_i} \ell(\sigma \epsilon_i; \theta_i) = h_i (\theta_i - \sigma \epsilon_i), \epsilon_i \sim \mathcal{N}(0, 1). \quad (13)$$

To validate the effectiveness of MaxVA, we compare MADAM with ADAM under a variety of different curvatures h and noise level σ on a 2D NQM (setting $d = 2$). For each setting of h and σ , we test both algorithms on a variety of learning rates. For ADAM, we additionally choose the best β and report the best results. See Appendix F for details. We run each setting 100 times to report the mean and standard error. MADAM consistently achieves 30-40% lower average loss with smaller standard error in all settings. Figure 2 shows the results for one of the settings, from which we find the best result of MADAM is better than ADAM under any choice of β , confirming the advantage of choosing an adaptive β_t . From the qualitative example, MaxVA also demonstrates smaller variance near convergence, enabled by a quicker response to impede the noise with a smaller β_t . More experimental results under other settings are provided in Appendix F.

5 EXPERIMENTS ON PRACTICAL DATASETS

In this section, we thoroughly evaluate MADAM and LAMADAM on a variety of tasks against well-calibrated baselines: CIFAR10/100 and ImageNet for image classification, IWSLT’14 DE-EN/WMT’16 EN-DE for neural machine translation, and the GLUE benchmark for natural language

Model	CIFAR-10	CIFAR-100	ImageNet
SGD	95.44 (.04)	79.62 (.07)	70.18
ADAM	95.37 (.03)	78.77 (.07)	66.54
LAPROP	95.34 (.03)	78.36 (.07)	70.02
MADAM	95.51 (.09)	79.32 (.08)	69.96
LAMADAM	95.38 (.11)	79.21 (.11)	70.16

Table 1: Accuracies on CIFAR10/100 and ImageNet. CIFAR10/100 experiments are the median (standard error) over 4 runs.

Method	IWSLT’14 DE-EN	WMT’16 EN-DE
LAPROP	35.98 (0.06)	27.02
LAMADAM	36.09 (0.04)	27.11

Table 2: BLEU score of LAPROP and LAMADAM for training transformers on machine translation datasets. We report the median and standard error for IWSLT’14 over 5 runs.

understanding. In both computer vision and NLP tasks, after careful tunings, we find the decoupled weight decay (Loshchilov & Hutter, 2018) gives much better results for ADAM, MADAM, LAPROP and LAMADAM. Therefore, we use this approach in all our experiments. Across all the plots in this section, we define the average step size at time t as the average of $|\eta_t m_t / (\sqrt{v_t} + \epsilon)|$ for ADAM/MADAM and $|\eta_t m_t|$ for LAPROP/LAMADAM over all the entries.

5.1 IMAGE CLASSIFICATION

To evaluate the effectiveness of MaxVA for image classification, we compare with SGD, ADAM and LAPROP in training ResNet18 (He et al., 2016) on CIFAR10, CIFAR100 and ImageNet. On all the datasets, we perform a grid search for the learning rate and weight decay, and report the best results for each method in Table 1. For CIFAR10/100, we train ResNet18 with a batch size of 128 for 200 epochs. We also find AMSGrad (Reddi et al., 2018) improves the classification accuracy of all adaptive methods evaluated on CIFAR10/100, so we apply AMSGrad in all experiments with adaptive methods. On ImageNet, we use the implementation from torchvision and the default multi-step learning rate schedule. We do not use AMSGrad in this case. Further details are in Appendix G.

Despite achieving a marginal improvement on CIFAR10, adaptive methods often underperforms carefully tuned SGD on CIFAR100 and ImageNet when training popular architectures such as ResNet, as confirmed by Wilson et al. (2017); Zhang et al. (2019c); Liu et al. (2020). Nevertheless, with the proposed MaxVA, we shrink the gap between adaptive methods and carefully tuned SGD on these image classification datasets, and achieve top-1 accuracy very close to SGD on ImageNet. Note our results with ResNet18 is better than the recent AdaBelief’s results with ResNet34 on CIFAR10/CIFAR100 (95.51/79.32 vs. 95.30/77.30 approximately), as well as AdaBelief with ResNet18 on ImageNet (70.16 vs. 70.08) (Zhuang et al., 2020).

5.2 NEURAL MACHINE TRANSLATION

We train Transformers from scratch with LAPROP and LAMADAM on IWSLT’14 German-to-English (DE-EN) translation (Cettolo et al., 2014) and WMT’16 English-to-German (EN-DE) translation, based on the implementation of fairseq.¹ We do not compare with SGD, since it is unstable for Transformers (Zhang et al., 2019b). We also show in Appendix D that ADABOUND cannot achieve any good result without degenerating into ADAM. More implementation details are in Appendix H.

IWSLT’14 DE-EN has 160k training examples, on which we use a smaller Transformer with 512-dimensional word embeddings and 1024 FFN dimensions. We train it for 60k iterations, during which each batch has up to 4096 tokens. Results are listed in Table 2. Note the baseline’s BLEU score is already 1.22 higher than the best results reported in Liu et al. (2020) using the same model. Although LAMADAM is only marginally better than LAPROP, LAMADAM uses much smaller update size than LAPROP, and it is not able for LAPROP to achieve better results when we scale its learning rate to get similar update sizes as LAMADAM, as shown in Appendix H.

WMT’16 EN-DE has 4.5M training examples, where same as Ott et al. (2018), we use a larger Transformer with 1024-dimensional word embeddings and 4096 FFN dimensions. Each batch has up to 480k tokens. We train for 32k iterations using the same inverse square root learning rate schedule as Vaswani et al. (2017). We evaluate the *single model* BLEU on newstest2013, unlike Liu et al. (2020) where models in the last 20 epochs are averaged to get the results. As shown in Table 2, LAMADAM also achieves better results.

¹<https://github.com/pytorch/fairseq>

Method	MNLI (Acc)	QNLI (Acc)	QQP (Acc)	RTE (Acc)	SST-2 (Acc)	MRPC (Acc)	CoLA (Mcc)	STS-B (Pearson)
Reported	87.6	92.8	91.9	78.7	94.8	90.2	63.6	91.2
ADAM	87.70 (.03)	92.85 (.06)	91.80 (.03)	79.25 (.71)	94.75 (.08)	88.50 (.24)	61.92 (1.1)	91.17 (.13)
LAPROP	87.80 (.04)	92.85 (.13)	91.80 (.03)	78.00 (.46)	94.65 (.11)	89.20 (.20)	63.01 (.61)	91.17 (.06)
MADAM	87.90 (.08)	92.95 (.07)	91.85 (.03)	79.60 (.66)	94.85 (.12)	89.70 (.17)	63.33 (.60)	91.28 (.03)
LAMADAM	87.80 (.03)	93.05 (.05)	91.85 (.05)	80.15 (.64)	95.15 (.15)	90.20 (.20)	63.84 (.85)	91.36 (.04)

Table 3: Results (median and variance) on the dev sets of GLUE based on finetuning the RoBERTa-base model (Liu et al. (2019)), from 4 runs with the same hyperparameter but different random seeds.

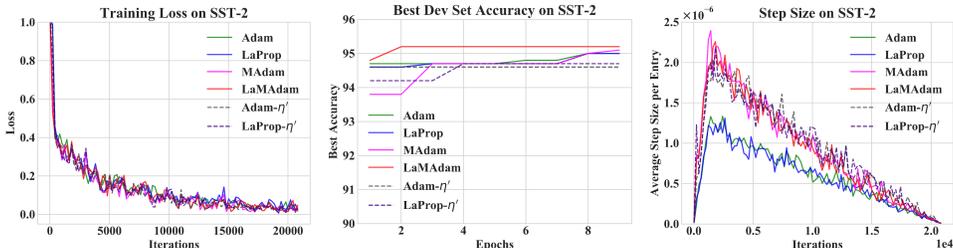


Figure 3: Training loss, validation accuracy and step size of various optimization methods on SST-2. All optimizers here use $\lambda = 0.1$. ADAM and LAPROP use $(\eta, \beta)=(1e-5, 0.98)$, MADAM and LAMADAM use $(\eta, \beta)=(4e-5, 0.5, 0.98)$, ADAM- η' and LAPROP- η' use $(\eta, \beta)=(1.6e-5, 0.98)$.

5.3 GENERAL LANGUAGE UNDERSTANDING EVALUATION (GLUE)

To evaluate MaxVA for transfer learning, we fine-tune pre-trained RoBERTa-base model (Liu et al., 2019) on 8 of the 9 tasks of the GLEU benchmark (Wang et al., 2018). Following prevalent validation settings (Devlin et al., 2019; Lan et al., 2020; Raffel et al., 2019), we report the median and standard error for fine-tuning the RoBERTa-base model (Liu et al., 2019) over 4 runs where only the random seeds are changed. The results are in Table 3. MADAM and LAMADAM give better scores than the corresponding baselines in the 8 tasks. More experimental details are in Appendix I.

To highlight the difference of the optimizers, we compare the training loss, dev set accuracy and the average step size on SST-2, as shown in Figure 3. Different from Machine Translation experiments where we train the Transformers from scratch, the adaptive step size of MADAM/LAMADAM is higher in this transfer learning setting. The ratio of the learning rate and step size of MaxVA to non-MaxVA optimizers are 4 and 1.8 respectively on GLUE, while on IWSLT’14 the two ratios are 2 and (approximately) 0.875. Because we start from a pre-trained model, the heavy tail of the gradient is alleviated, just as the BERT model in the later stage of training as shown by Zhang et al. (2019b), and the curvature of the loss landscape should be smaller. Therefore, MaxVA selects larger adaptive step sizes for better convergence. Same as in the Machine Translation experiments, the highest test accuracy of ADAM/LAPROP cannot reach the same value as MADAM/LAMADAM by simply scaling the base learning rate η to reach similar step sizes as MADAM/LAMADAM.

5.4 LARGE-BATCH PRETRAINING FOR BERT

We use the NVIDIA BERT pretraining repository to perform large-batch pretraining for BERT-Base model on the Wikipedia Corpus only.² Each run takes about 52 hours on 8 V100 GPUs. Training is divided into two phases: the first phase uses a batch size of 64K with input sequence length 128 for 7,038 steps; the second phase uses a batch size 32K with input sequence length 512 for 1563 steps. The total of steps is significantly smaller than the 1,000,000 steps used in the small-batch training of Devlin et al. (2019). Therefore, a faster adaptation to curvature in each step is more important.

This point is validated by the faster convergence of MADAM in both phases, as shown in the training loss curves in Figure 4. Contrary to the observation by You et al. (2020), ADAM even converges

²<https://github.com/NVIDIA/DeepLearningExamples/tree/master/PyTorch/LanguageModeling/BERT>.

Note the results from the repository are for BERT-Large trained with additional data from BookCorpus. We did not include BookCorpus due to the difficulty of obtaining the data.

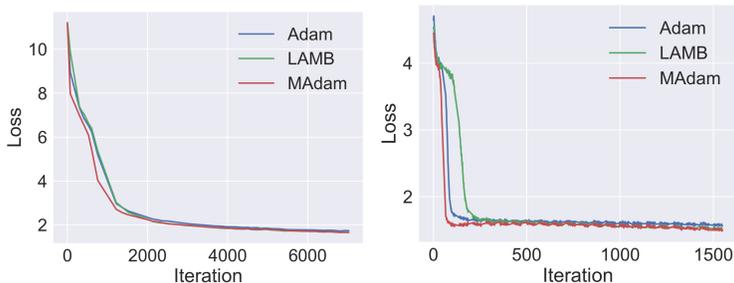


Figure 4: Training losses of ADAM, LAMB and MADAM on Wikipedia Corpus in the two training phases.

faster than LAMB in the earlier iterations. You et al. (2020) only explored weight decay of up to 0.01 for ADAM, but we find using larger weight decay of 0.1 together with gradient clipping ($\|g_t\|_2 \leq 1$, same as LAMB) stabilizes ADAM in this large-batch setting. So we use a weight decay of 0.1 and the same gradient clipping for both MADAM and ADAM. For MADAM and ADAM, we do a grid search on the learning rate of phase 1 while keeping the ratios of learning rate in phase 1 and phase 2 to the same as LAMB. We use $\bar{\beta} = 0.999$, $\underline{\beta} = 0.5$ for MADAM. For LAMB, we use the default setting from the aforementioned repository.

The faster adaptation of MaxVA improves the stability, which enables MADAM to use a much larger learning rate to achieve faster convergence than ADAM. The best learning rate for MADAM is $3.4e-3$. We tried learning rates in $\{7e-4, 8e-4, 9e-4, 1e-3\}$ for ADAM, and find it always diverges when the learning rate is higher or equal to $9e-4$. The best result of ADAM is achieved with learning rate $8e-4$. MADAM achieves a training loss of 1.492, while LAMB achieves a training loss of 1.507, and ADAM has the worst training loss 1.568. The test scores of the models pretrained with MADAM/LAMB/ADAM are 88.53/87.60/88.07 (F1) and 82.10/81.40/80.78 (Accuracy) on SQuAD v1.1 and MNLI, respectively.

6 RELATED WORK

Various adaptive methods have been proposed and broadly applied in deep learning (Kingma & Ba, 2015; Duchi et al., 2011; Tieleman & Hinton, 2012; Zeiler, 2012). Reddi et al. (2018) proposed to compute the adaptive learning rate with the coordinate-wise maximum value of the running squared gradient so that the adaptive learning rate does not increase. ADABOUND (Luo et al., 2019) clips the adaptive learning rate of ADAM with a decreasing upper bound and an increasing lower bound. Lookahead (Zhang et al., 2019c) computes weight updates by looking ahead at the sequence of “fast weights” generated by another optimizer. Padam (Chen et al., 2018) improves the generalization of adaptive methods by choosing a proper exponent for the v_t of AMSGRAD. LAPROP (Ziyin et al., 2020) uses local running estimation of the variance to normalize the gradients, resulting in higher empirical stability. You et al. (2017) proposes Layer-wise Adaptive Rate Scaling (LARS), and scales the batch size to 16,384 for training ResNet50. LAMB (You et al., 2020) applies a similar layer-wise learning rate on ADAM to improve LARS on training BERT. Starting from a similar motivation of adapting to the curvature, the recent work AdaBelief (Zhuang et al., 2020) directly estimates the exponential running average of the gradient deviation to compute the adaptive step sizes. Our approach finds the averaging coefficients β_t automatically by maximizing the estimated variance for a faster adaptation to the curvature, which could be complementary to all the aforementioned methods, and is the first to explore in this direction to our knowledge.

7 CONCLUSION

In this paper, we present Maximum Variation Averaging (MaxVA), a novel adaptive learning rate scheme that replaces the exponential running average of squared gradient with an adaptive weighted mean. In each step, MaxVA chooses the weight β_t for each coordinate, such that the estimated gradient variance is maximized. This enables MaxVA to: (1) take smaller steps when large curvatures or abnormal gradients are present, which leads to more desirable convergence behaviors in face of noisy gradients; (2) adapt faster to the geometry of the objective, achieving faster convergence in the large-batch setting. We illustrate how our method improves convergence by a better adaptation to variance, and demonstrate strong empirical results on a wide range of tasks. We prove MaxVA converges in the nonconvex stochastic optimization setting under mild assumptions.

REFERENCES

- Eneko Agirre, Lluís M´arquez, and Richard Wicentowski (eds.). *Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)*. Association for Computational Linguistics, 2007.
- Lukas Balles and Philipp Hennig. Dissecting adam: The sign, magnitude and variance of stochastic gradients. In *ICML*, pp. 404–413, 2018.
- Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. The fifth PASCAL recognizing textual entailment challenge. In *TAC*, 2009.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *ICML*, pp. 560–569, 2018.
- Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. Report on the 11th iwslt evaluation campaign, iwslt 2014. In *IWSLT*, volume 57, 2014.
- Jinghui Chen, Dongruo Zhou, Yiqi Tang, Ziyang Yang, and Quanquan Gu. Closing the generalization gap of adaptive gradient methods in training deep neural networks. *arXiv preprint arXiv:1806.06763*, 2018.
- Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of adam-type algorithms for non-convex optimization. *ICLR*, 2019.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. The PASCAL recognising textual entailment challenge. In *Machine learning challenges. evaluating predictive uncertainty, visual object classification, and recognising textual entailment*. Springer, 2006.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, pp. 4171–4186, 2019.
- William B Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the International Workshop on Paraphrasing*, 2005.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, pp. 2121–2159, 2011.
- Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third PASCAL recognizing textual entailment challenge. In *Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing*, 2007.
- R Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szepesky. The second pascal recognising textual entailment challenge. In *Proceedings of the Second PASCAL Challenges Workshop on Recognising Textual Entailment*, 2006.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pp. 770–778, 2016.
- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- Shankar Iyer, Nikhil Dandekar, and Kornel Csernai. First quora dataset release: Question pairs, 2017. URL <https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs>.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun (eds.), *ICLR*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *ICLR*, 2020.

- Hector J Levesque, Ernest Davis, and Leora Morgenstern. The Winograd schema challenge. In *AAAI Spring Symposium: Logical Formalizations of Commonsense Reasoning*, 2011.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. *ICLR*, 2020.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv:1907.11692*, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2018.
- Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic bound of learning rate. *ICLR*, 2019.
- Jerry Ma and Denis Yarats. On the adequacy of untuned warmup for adaptive optimization. *arXiv:1910.04209*, 2019.
- Myle Ott, Sergey Edunov, David Grangier, and Michael Auli. Scaling neural machine translation. In *WMT*, pp. 1–9, 2018.
- Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le. SpecAugment: A simple data augmentation method for automatic speech recognition. *Interspeech*, pp. 2613–2617, 2019.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv:1910.10683*, 2019.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *EMNLP*, 2016.
- Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *ICLR*, 2018.
- Tom Schaul, Sixin Zhang, and Yann LeCun. No more pesky learning rates. In *ICML*, pp. 343–351, 2013.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*, 2013.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- T. Tieleman and G. Hinton. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *EMNLP*, pp. 353, 2018.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. Neural network acceptability judgments. *arXiv preprint 1805.12471*, 2018.
- Adina Williams, Nikita Nangia, and Samuel R. Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *NAACL*, 2018.
- Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. In *Neurips*, pp. 4148–4158, 2017.

- Yuhuai Wu, Mengye Ren, Renjie Liao, and Roger Grosse. Understanding short-horizon bias in stochastic meta-optimization. *arXiv:1803.02021*, 2018.
- Yang You, Igor Gitman, and Boris Ginsburg. Scaling SGD batch size to 32k for imagenet training. *CoRR*, abs/1708.03888, 2017.
- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. In *ICLR*, 2020.
- Manzil Zaheer, Sashank Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. Adaptive methods for nonconvex optimization. In *NeurIPS*, pp. 9793–9803, 2018.
- Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.
- Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George Dahl, Chris Shallue, and Roger B Grosse. Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model. In *NeurIPS*, pp. 8194–8205, 2019a.
- Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank J Reddi, Sanjiv Kumar, and Suvrit Sra. Why adam beats sgd for attention models. *arXiv:1912.03194*, 2019b.
- Michael R. Zhang, James Lucas, Jimmy Ba, and Geoffrey E. Hinton. Lookahead optimizer: k steps forward, 1 step back. In *NeurIPS*, pp. 9593–9604, 2019c.
- Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. Freelib: Enhanced adversarial training for natural language understanding. In *ICLR*, 2020.
- Juntang Zhuang, Tommy Tang, Sekhar Tatikonda, Nicha Dvornek, Yifan Ding, Xenophon Papademetris, and James S. Duncan. Adabelief optimizer: Adapting stepsizes by the belief in observed gradients. In *NeurIPS*, 2020.
- Liu Ziyin, Zhikang T Wang, and Masahito Ueda. Laprop: a better way to combine momentum with adaptive gradient. *arXiv:2002.04839*, 2020.

ADAPTIVE LEARNING RATES WITH MAXIMUM VARIATION AVERAGING (APPENDIX)

A DERIVING THE CLOSED FORM SOLUTION EQ.7

Plugging Eq. 3,4,5, and the unbiased estimations $u_t(\beta) = \tilde{u}_t(\beta)/w_t(\beta)$, $v_t(\beta) = \tilde{v}_t(\beta)/w_t(\beta)$ into Eq. 6, each coordinate is solving the same problem:

$$\arg \max_{\beta} f(\beta) = \frac{\beta w_{t-1} v_{t-1} + (1-\beta) g_t^2}{\beta w_{t-1} + (1-\beta)} - \left[\frac{\beta w_{t-1} u_{t-1} + (1-\beta) g_t}{\beta w_{t-1} + (1-\beta)} \right]^2. \quad (14)$$

Let $\gamma = 1/[\beta w_{t-1} + (1-\beta)] \in [1, 1/w_{t-1}]$, we can see $f(\beta)$ can be represented as a quadratic function of γ . Specifically,

$$f(\beta) = h(\gamma) = \frac{w_{t-1} v_{t-1} - g_t^2}{w_{t-1} - 1} + \left[g_t^2 - \frac{w_{t-1} v_{t-1} - g_t^2}{w_{t-1} - 1} \right] \gamma - \left\{ \frac{w_{t-1} u_{t-1} - g_t}{w_{t-1} - 1} + \left[g_t - \frac{w_{t-1} u_{t-1} - g_t}{w_{t-1} - 1} \right] \gamma \right\}^2.$$

Meanwhile, β is a monotonic function of γ . Therefore, $f(\beta)$ has a unique maximum value.

To find the maximum value, we return to Eq. 14, from which we can find a stationary point

$$\frac{v_{t-1} - u_{t-1}^2 + (g_t - u_{t-1})^2}{w_{t-1} [(g_t - u_{t-1})^2 - v_{t-1} + u_{t-1}^2] + v_{t-1} - u_{t-1}^2 + (g_t - u_{t-1})^2}. \quad (15)$$

B CONVERGENCE PROOF

Following the convergence proofs of YOGI (Zaheer et al., 2018), we prove the convergence of MADAM in the nonconvex setting.

Proof of Theorem 1.

Proof. Recall that we have assumed the update steps of MADAM as

$$\theta_{t+1,i} = \theta_{t,i} - \eta_t \frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon}, \quad (16)$$

for all $i \in [d]$, and that f is L -smooth, which results in the following inequalities:

$$\begin{aligned} f(\theta_{t+1}) &\leq f(\theta_t) + \langle \nabla f(\theta_t), \theta_{t+1} - \theta_t \rangle + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \\ &= f(\theta_t) - \eta_t \sum_{i=1}^d \nabla f(\theta_{t,i}) \frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon} + \frac{L\eta_t^2}{2} \sum_{i=1}^d \frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2}. \end{aligned} \quad (17)$$

Note the stochastic gradient is defined as $g_t = \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla_{\theta} \ell(x_j; \theta_t)$. Given θ_t , we take expectation over the stochastic gradient g_t in Eq. 17 (denoted as $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot|\theta_t]$) to get

$$\begin{aligned}
\mathbb{E}_t[f(\theta_{t+1})] &\leq f(\theta_t) - \eta_t \sum_{i=1}^d \left([\nabla f(\theta_t)]_i \times \mathbb{E}_t \left[\frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon} \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2} \right] \\
&= f(\theta_t) - \eta_t \sum_{i=1}^d \left([\nabla f(\theta_t)]_i \times \mathbb{E}_t \left[\frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon} - \frac{g_{t,i}}{\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon} + \frac{g_{t,i}}{\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon} \right] \right) \\
&\quad + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2} \right] \\
&= f(\theta_t) - \eta_t \sum_{i=1}^d \left([\nabla f(\theta_t)]_i \times \left[\mathbb{E}_t \left[\frac{g_{t,i}}{\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon} \right] + \mathbb{E}_t \left[\frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon} - \frac{g_{t,i}}{\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon} \right] \right] \right) \\
&\quad + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2} \right] \\
&\leq f(\theta_t) - \eta_t \sum_{i=1}^d \left[\frac{[\nabla f(\theta_t)]_i^2}{\sqrt{\beta}v_{t-1,i} + \epsilon} - \frac{\sigma G}{\epsilon\sqrt{t}} \right] + \eta_t \sum_{i=1}^d [\nabla f(\theta_t)]_i \mathbb{E}_t \underbrace{\left[\frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon} - \frac{g_{t,i}}{\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon} \right]}_{T_1} \\
&\quad + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2} \right], \tag{18}
\end{aligned}$$

where the second equality holds by applying Lemma 2 to the first expectation term, and taking the absolute value of the second expectation term.

Next, we need to bound the term T_1 to show convergence. First, we have the following upper bound for T_1 :

$$\begin{aligned}
T_1 &= \left| \frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon} - \frac{g_{t,i}}{\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon} \right| \\
&\leq |g_{t,i}| \times \left| \frac{1}{\sqrt{v_{t,i}} + \epsilon} - \frac{1}{\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon} \right| \\
&\leq \frac{|g_{t,i}|}{(\sqrt{v_{t,i}} + \epsilon)(\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon)} \times \left| \frac{v_{t,i} - \beta_{t,i}v_{t-1,i}}{\sqrt{v_{t,i}} + \sqrt{\beta_{t,i}v_{t-1,i}}} \right| \\
&= \frac{|g_{t,i}|}{(\sqrt{v_{t,i}} + \epsilon)(\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon)} \times \frac{(1 - \beta_{t,i})g_{t,i}^2}{\sqrt{v_{t,i}} + \sqrt{\beta_{t,i}v_{t-1,i}}}, \tag{19}
\end{aligned}$$

where the last equality comes from the definition of $v_{t,i} = \beta_{t,i}v_{t-1,i} + (1 - \beta_{t,i})g_{t,i}^2$. We can further bound T_1 as

$$\begin{aligned}
T_1 &= \frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)(\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon)} \times \frac{(1 - \beta_{t,i})|g_{t,i}|}{\sqrt{\beta_{t,i}v_{t-1,i} + (1 - \beta_{t,i})g_{t,i}^2} + \sqrt{\beta_{t,i}v_{t-1,i}}} \\
&\leq \frac{g_{t,i}^2}{(\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon)\epsilon} \times \frac{(1 - \beta_{t,i})|g_{t,i}|}{\sqrt{(1 - \beta_{t,i})g_{t,i}^2}} \\
&= \frac{\sqrt{1 - \beta_{t,i}}g_{t,i}^2}{(\sqrt{\beta_{t,i}v_{t-1,i}} + \epsilon)\epsilon} \leq \frac{\sqrt{1 - \underline{\beta}}g_{t,i}^2}{(\sqrt{\underline{\beta}v_{t-1,i}} + \epsilon)\epsilon} \tag{20}
\end{aligned}$$

Since the loss on each sample s satisfies $|\nabla \ell(x, s)|_i \leq G$, we will have $|\nabla f(x)|_i \leq G$ for $\forall i \in [d]$. Substituting the coefficients of T_1 in Eq. 18 with this gradient bound, we have

$$\begin{aligned}
\mathbb{E}_t[f(\theta_{t+1})] &\leq f(\theta_t) - \eta_t \sum_{i=1}^d \frac{[\nabla f(\theta_t)]_i^2}{\sqrt{\beta}v_{t-1,i} + \epsilon} + \frac{\eta_t G \sqrt{1-\underline{\beta}}}{\epsilon} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{\sqrt{\beta}v_{t-1,i} + \epsilon} \right] \\
&\quad + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2} \right] + \frac{\sigma\eta dG}{\epsilon\sqrt{t}}, \\
&\leq f(\theta_t) - \eta_t \sum_{i=1}^d \frac{[\nabla f(\theta_t)]_i^2}{\sqrt{\beta}v_{t-1,i} + \epsilon} + \frac{\eta_t G \sqrt{1-\underline{\beta}}}{\epsilon} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{\sqrt{\beta}v_{t-1,i} + \epsilon} \right] \\
&\quad + \frac{L\eta_t^2}{2\epsilon} \sum_{i=1}^d \mathbb{E}_t \left[\frac{g_{t,i}^2}{\sqrt{\beta}v_{t-1,i} + \epsilon} \right] + \frac{\sigma\eta dG}{\epsilon\sqrt{t}}, \\
&\leq f(\theta_t) + \underbrace{\sum_{i=1}^d \left(-\frac{\eta_t}{(\sqrt{\beta}v_{t-1,i} + \epsilon)} + \frac{\eta_t G \sqrt{1-\underline{\beta}}}{\epsilon(\sqrt{\beta}v_{t-1,i} + \epsilon)} + \frac{L\eta_t^2}{2\epsilon(\sqrt{\beta}v_{t-1,i} + \epsilon)} \right)}_{T_2} [\nabla f(\theta_t)]_i^2 \\
&\quad + \left(\frac{\eta_t G \sqrt{1-\underline{\beta}}}{\epsilon} + \frac{L\eta_t^2}{2\epsilon} \right) \sum_{i=1}^d \frac{\sigma^2}{b_t(\sqrt{\beta}v_{t-1,i} + \epsilon)} + \frac{\sigma\eta dG}{\epsilon\sqrt{t}}, \tag{21}
\end{aligned}$$

where the second inequality comes from the fact that $\sqrt{v_{t,i}} + \epsilon \geq \epsilon$ and $v_{t,i} = \beta_{t,i}v_{t-1,i} + (1 - \beta_{t,i})g_{t,i}^2 \geq \beta v_{t-1,i}$, and the third inequality comes from applying Lemma 1 by Zaheer et al. (2018) to $\mathbb{E}_t[g_t^2]$. The application of Lemma 1 is possible because v_{t-1} is independent of the t -th batch. By the assumptions for $\epsilon, G, \underline{\beta}$, we have

$$\frac{L\eta_t^2}{2\epsilon} \leq \frac{\eta}{4}, \quad \frac{\eta_t G \sqrt{1-\underline{\beta}}}{\epsilon} \leq \frac{1}{4}\eta. \tag{22}$$

Plugging these two results and the assumption $\bar{\beta} \leq 2\underline{\beta}$ into T_2 , we have

$$\begin{aligned}
T_2 &\leq -\frac{\eta}{\sqrt{\beta}v_{t-1,i} + \epsilon} + \frac{\eta}{2(\sqrt{\beta}v_{t-1,i} + \epsilon)} \\
&\leq -\frac{\eta}{\sqrt{2}(\sqrt{\beta}v_{t-1,i} + \epsilon)} + \frac{\eta}{2(\sqrt{\beta}v_{t-1,i} + \epsilon)} \\
&\leq \frac{\eta}{5(\sqrt{\beta}v_{t-1,i} + \epsilon)}
\end{aligned} \tag{23}$$

the main inequality, we have

$$\begin{aligned}
\mathbb{E}_t[f(\theta_{t+1})] &\leq f(\theta_t) - \frac{\eta}{5} \sum_{i=1}^d \frac{[\nabla f(\theta_t)]_i^2}{\sqrt{\beta}v_{t-1,i} + \epsilon} + \left(\frac{\eta_t G \sqrt{1-\underline{\beta}}}{\epsilon} + \frac{L\eta_t^2}{2\epsilon} \right) \sum_{i=1}^d \frac{\sigma^2}{t(\sqrt{\beta}v_{t-1,i} + \epsilon)} + \frac{\sigma\eta dG}{\epsilon\sqrt{t}} \\
&\leq f(\theta_t) - \frac{\eta}{5(G\sqrt{\underline{\beta}} + \epsilon)} \|\nabla f(\theta_t)\|^2 + \left(\frac{\eta_t G \sqrt{1-\underline{\beta}}}{\epsilon^2} + \frac{L\eta_t^2}{2\epsilon^2} \right) \frac{\sigma^2 d}{t} + \frac{\sigma\eta dG}{\epsilon\sqrt{t}}, \tag{24}
\end{aligned}$$

where we have replaced b_t with t by our assumption on the batch size, and the second inequality comes from the fact that $v_{t-1,i} \leq G^2$. Taking expectation on both the LHS and RHS for the inequalities at $t = 1, \dots, T$, using telescope sum and rearranging the terms, we can conclude that

$$\frac{\eta}{5(G\sqrt{\underline{\beta}} + \epsilon)} \sum_{i=1}^T \|\nabla f(\theta_t)\|^2 \leq f(\theta_1) - E[f(\theta_{T+1})] + \left(\frac{\eta G \sqrt{1-\underline{\beta}}}{\epsilon^2} + \frac{L\eta^2}{2\epsilon^2} \right) \sigma^2 d \log(T+1) + \frac{2\sigma\eta dG}{\epsilon} \sqrt{T}. \tag{25}$$

Multiplying both sides with $\frac{5(G\sqrt{\beta+\epsilon})}{T\eta}$, and using the fact that $f(x^*) \leq f(\theta_{t+1})$, we conclude that

$$\frac{1}{T} \sum_{i=1}^T \|\nabla f(\theta_t)\|_i^2 \leq 5(G\sqrt{\beta+\epsilon}) \left(\frac{f(\theta_1) - f(x^*)}{\eta T} + \left(\frac{G\sqrt{1-\beta}}{\epsilon^2} + \frac{L\eta}{2\epsilon^2} \right) \frac{\sigma^2 d \log(T+1)}{T} + \frac{2\sigma d G}{\epsilon\sqrt{T}} \right). \quad (26)$$

Lemma 2. Assume the gradient is bounded as $\|\nabla_{\theta} \ell(x; \theta)\|_{\infty} \leq G$, and has bounded variance $\mathbb{E}[(\nabla_{\theta} \ell(x; \theta))_i - [\nabla f(\theta_t)]_i]^2 \leq \sigma^2$, and the batch size $b_t = t$. For the t -th iteration of MADAM, we have

$$-[\nabla f(\theta_t)]_i \mathbb{E}_t \left[\frac{g_{t,i}}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right] \leq -\frac{[\nabla f(\theta_t)]_i^2}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} + \frac{G\sigma}{\epsilon\sqrt{t}} \quad (27)$$

Proof. The LHS can be decomposed as

$$\begin{aligned} \text{LHS} &= -[\nabla f(\theta_t)]_i \mathbb{E}_t \left[\frac{[\nabla f(\theta_t)]_i}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right] - [\nabla f(\theta_t)]_i \mathbb{E}_t \left[\frac{g_{t,i} - [\nabla f(\theta_t)]_i}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right] \\ &\leq -\frac{[\nabla f(\theta_t)]_i^2}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} - \underbrace{[\nabla f(\theta_t)]_i \mathbb{E}_t \left[\frac{g_{t,i} - [\nabla f(\theta_t)]_i}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right]}_{T_3}, \end{aligned} \quad (28)$$

where the inequality comes from taking the upper bound of $\beta_{t,i}$, since the first term is non-positive. Let $[h(x)]_+$ and $[h(x)]_-$ be the operators for taking the positive and negative values of function $h(x)$ respectively, i.e.,

$$[h(x)]_+ = \begin{cases} h(x), & \text{if } h(x) > 0 \\ 0, & \text{otherwise} \end{cases}, \quad [h(x)]_- = \begin{cases} h(x), & \text{if } h(x) < 0 \\ 0, & \text{otherwise} \end{cases}. \quad (29)$$

It is obvious that $\mathbb{E}[[X]_+] \leq \mathbb{E}[|X|] \leq \sqrt{\mathbb{E}[X^2]}$, where the second inequality comes from Cauchy-Schwarz inequality. Similarly, $\mathbb{E}[[X]_-] \geq -\mathbb{E}[|X|] \geq -\sqrt{\mathbb{E}[X^2]}$. With this in mind, we have

$$0 \leq \mathbb{E}_t \left[[g_{t,i} - [\nabla f(\theta_t)]_i]_+ \right] \leq \sqrt{\mathbb{E}_t [g_{t,i} - [\nabla f(\theta_t)]_i]^2} \leq \frac{\sigma}{\sqrt{t}}, \quad (30)$$

where the last inequality comes from applying Lemma 1 from Zaheer et al. (2018) under the bounded gradient variance assumption, and the assumption that the batch size grows as $b_t = t$. Similarly, we have

$$-\frac{\sigma}{\sqrt{t}} \leq \mathbb{E}_t \left[[g_{t,i} - [\nabla f(\theta_t)]_i]_- \right] \leq 0. \quad (31)$$

Now we will decompose and bound T_3 as

$$\begin{aligned} T_3 &= -\mathbb{E}_t \left[[\nabla f(\theta_t)]_i \frac{[g_{t,i} - [\nabla f(\theta_t)]_i]_+}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right] - \mathbb{E}_t \left[[\nabla f(\theta_t)]_i \frac{[g_{t,i} - [\nabla f(\theta_t)]_i]_-}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right] \\ &\leq \begin{cases} -\mathbb{E}_t \left[[\nabla f(\theta_t)]_i \frac{[g_{t,i} - [\nabla f(\theta_t)]_i]_-}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right], & \text{if } [\nabla f(\theta_t)]_i > 0 \\ -\mathbb{E}_t \left[[\nabla f(\theta_t)]_i \frac{[g_{t,i} - [\nabla f(\theta_t)]_i]_+}{\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon}} \right], & \text{otherwise} \end{cases} \\ &\leq \frac{\sigma |\nabla f(\theta_t)_i|}{(\sqrt{\beta_{t,i} v_{t-1,i} + \epsilon})\sqrt{t}}, \\ &\leq \frac{\sigma G}{\epsilon\sqrt{t}}. \end{aligned} \quad (32)$$

Plugging this inequality back into Eq. 28 and we will get the RHS. \square

C PRACTICAL NOTES OF β_t

Claims and arguments:

- For $t > 1$, since $0 < \beta_t \leq 1$, w_t will monotonically increase from $(1 - \beta_1)$ to 1.
This is obvious since in every step, w_t is an interpolation between w_{t-1} and 1, and $w_t \geq w_{t-1}$. We have also set $w_1 = 1 - \beta_1$.
- For any g_t, u_{t-1}, v_{t-1} satisfying $v_{t-1} - u_{t-1}^2 > 0$ in Eq. 7, we have $\beta_t \in [1/(1 + w_{t-1}), 1/(1 - w_{t-1})]$.
Eq. 9 is monotonic in R_t . Since g_t can be any value, R_t can be any value from 0 to ∞ . If $R_t = 0$, β_t takes the largest value $1/(1 - w_t)$. If $R_t \rightarrow \infty$, $\beta_t \rightarrow 1/(w_t + 1)$.
- As $t \rightarrow \infty$, $w_t \rightarrow 1$ and $\beta_t \in [0.5, \infty]$.
Combining Claims 1 and 2 to get this result.
- Adding a small coefficient $\delta > 0$ to the denominator of Eq. 7 has negligible effect on the value of β_t and does not violate the maximum variation objective (Eq. 6).
Since δ is small, it has negligible effect on β_t when division by zero does not happen. We only need to confirm adding δ will not affect the solution when division by zero happens. We can re-write the dividend of Eq. 7 as

$$(w_{t-1} + 1)(g_t - u_{t-1})^2 + (1 - w_{t-1})(v_{t-1} - u_{t-1}^2). \quad (33)$$

Since $\mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \text{Var}[X] \geq 0$, we can conclude that $v_{t-1} - u_{t-1}^2 \geq 0$.

When $1 - \beta_1 \leq w_{t-1} < 1$, Eq. 33 can be 0 only when $g_t = u_{t-1}$ and $v_{t-1} = u_{t-1}^2$. In this special case, we can set β_t to any value in $[0, 1]$ without changing $\tilde{\sigma}_t^2$; we will always have $v_t = \tilde{v}_{t-1}/w_{t-1} = v_{t-1}$, $u_t = \tilde{u}_{t-1}/w_{t-1} = u_{t-1}$, and $\tilde{\sigma}_t^2 = 0$. Only $w_t = (w_{t-1} - 1)\beta_t + 1$ is affected by β_t , which takes a larger value when β_t is smaller. The solution given by adding δ to the denominator is $\beta_t = 0$, and the following clipping will set $\beta_t = \bar{\beta}$, resulting in the largest possible $w_t = (w_{t-1} - 1)\bar{\beta} + 1$. In the next step, if Eq. 33 is not zero, then we have $\beta_{t+1} = 1/(w_t + 1)$, and we know $g_{t+1} \neq u_t$.³ In this case, for $0.5 \leq \beta_{t+1} < 1$, $\tilde{\sigma}_{t+1}^2$ increases as β_{t+1} decreases, so setting w_t to its maximum will achieve the maximum variance at the next step. Otherwise if Eq. 33 is zero, doing this will not change $\tilde{\sigma}_{t+1}^2 = 0$.

When $w_{t-1} = 1$, Eq. 33 is 0 if and only if $g_t = u_{t-1}$. As a result, if $v_{t-1} = u_{t-1}^2$, we have the same conclusion as before. Otherwise, $\beta_t = (v_{t-1} - u_{t-1}^2)/\delta$ before clipping, and $\beta_t = \bar{\beta}$ after clipping. Also, any $0 < \beta_t < 1$ will not change the value of $u_t = \beta_t u_{t-1} + (1 - \beta_t)g_t = u_{t-1}$. Since $g_t^2 = u_{t-1}^2 < v_{t-1}$, to maximize $\tilde{\sigma}_t^2 = v_t(\beta) - u_t^2$, we should set $\beta_t = \bar{\beta}$ so that $v_t(\beta)$ takes the maximum value, which is consistent with the solution after adding δ to the denominator.

D ADABOUND MIGHT FAIL ON TRANSFORMERS?

Since SGD often performs much worse than Adam on transformers, and ADABOUND transitions into SGD asymptotically, it is reasonable to believe that AdaBound would not converge well on transformers. We did experiments on the IWSLT'14 dataset to evaluate ADABOUND on Transformers. AdaBound clips the effective step size to be within $0.1 - \frac{0.1}{\gamma t + 1}$ and $0.1 - \frac{0.1}{\gamma t}$, and recommends setting $\gamma = 1 - \beta_2 = 10^{-3}$. In practice, this setting only gives a < 24 test BLEU on IWSLT'14. To explore the full potential of AdaBound, we tried $\gamma \in \{10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}\}$, and found $\gamma = 10^{-8}$ to give the best BLEU 35.99 (0.04). However, as shown in Figure 5, AdaBound does not effectively clip most of the coordinates even in the last iteration with $\gamma = 10^{-8}$, which means ADABOUND essentially degraded into Adam, yet it gives better results than those effectively doing clipping. By comparison, the best result of MAdam and Adam with AMSGrad is 36.07(0.07) / 35.87 (0.05), respectively.

³Otherwise we will still have $g_{t+1} = u_t$, $g_{t+1}^2 = u_t^2 = v_t$ and Eq. 33 is 0.

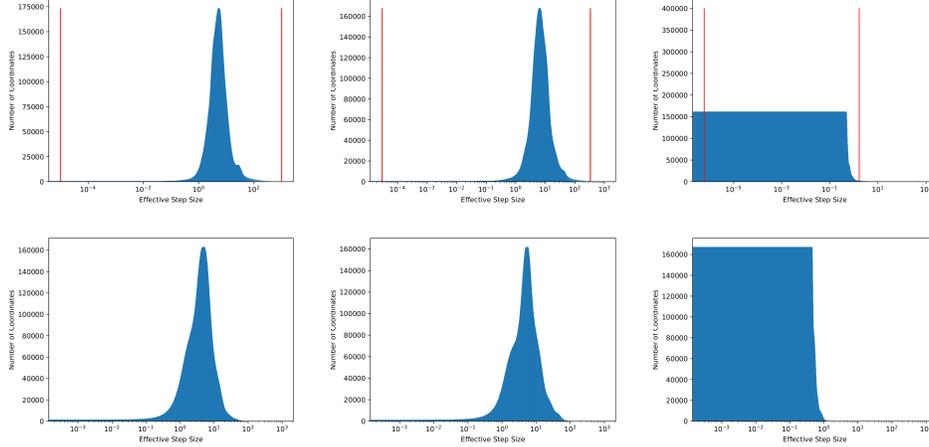


Figure 5: Distribution of effective step size of AdaBound and MADam at iteration 10000, 30000 and 60000 on IWSLT’14. Red lines indicate the clipping range of AdaBound. On the top/bottom are results of AdaBound/MADam with learning rates $5e-4/1.25e-3$.

E EXPERIMENTAL DETAILS ON THE SYNTHETIC FINITE-SAMPLE EXPERIMENT

Same as Chen et al. (2019), we use constant learning rates η in every step, and set $\alpha = 0, \beta = 0.9$ for ADAM and AMSGRAD. For MADAM, we set $\alpha = 0, (\beta, \bar{\beta}) = (0.5, 1)$. ADAM never converged for a variety of η we tried within $[10^{-4}, 1]$, consistent with Chen et al. (2019). Generally, a larger η gives faster convergence for both AMSGRAD and MADAM. For reproducibility, we repeat experiments 100 times with the same settings, and choose the η for AMSGRAD and MADAM where the solution $|\theta^*| < 0.1$ every time. $\eta = 1.2$ satisfies this requirement for MADAM, but AMSGRAD only satisfied it 1% of the times for $\eta = 1.2$ and 65% of the times for $\eta = 0.9$. $\eta = 0.8$ is the largest η we find for AMSGRAD to achieve 100% satisfaction. Therefore, we use $\eta = 0.8$ for both ADAM and AMSGRAD.

F DETAILS AND ADDITIONAL EXPERIMENTAL RESULTS ON THE NOISY QUADRATIC MODEL

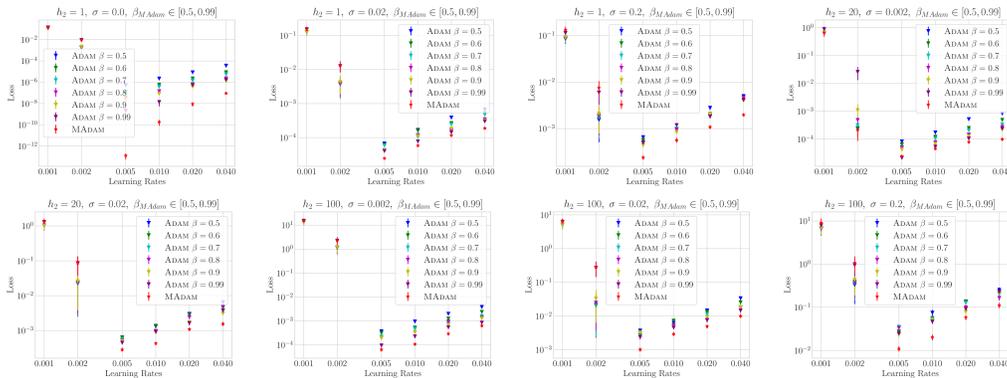


Figure 6: More results on the Noisy Quadratic Model.

Details of experimental settings We set $\beta = 0.5, \bar{\beta} = 0.99$ for MADAM, and for fair comparison, we do a grid search for ADAM with $\beta \in [0.5, 0.6, 0.7, 0.8, 0.9, 0.99]$, and only report the results with the best β . We repeat the experiments 100 times under each setting, where we select a random in-

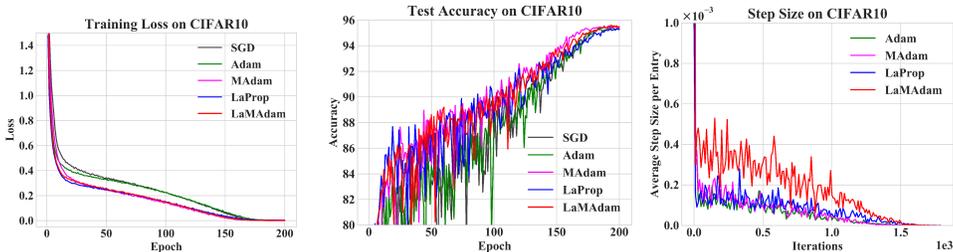


Figure 7: Training loss, test accuracy and average step size on CIFAR10.

tialization of $\theta \sim \mathcal{N}(0, I)$ each time, and run MADAM and ADAM with different hyper-parameters from this random initialization. Each run takes 1000 iterations by default.

Additional Results We give more results comparing ADAM and MADAM on the Noisy Quadratic Model. The results are shown in Figure 6. Generally, the best result of MADAM has a more significant margin when h_2 and σ are higher, i.e., the improvement is more significant when the problem is worse conditioned and the noise level is higher. Note that for each trial, we start both algorithms from the same random initialization.

G ADDITIONAL EXPERIMENTAL RESULTS AND DETAILS ON IMAGE CLASSIFICATION.

Additional Experimental Results In Figure 7, we plot the training loss, test accuracy and average step size on CIFAR10. We see that LAMADAM and MADAM produce slightly better curves than others.

Model, learning rate schedules and data augmentations On CIFAR10 and CIFAR100, the ResNet18 comes from a public repository,⁴ which has a base width of 64 by default. We use random cropping (4-pixel zero paddings on each side) and random horizontal flip as the data augmentations. Instead of using the multi-step schedule, we find the cosine learning rate schedule to yield better results for both SGD and adaptive methods. Therefore, we use the cosine learning rate schedule and set a final learning rate of $2e-6$ in all cases. On ImageNet, we use random resized crop and random horizontal flip for data augmentation. For the multi-step learning rate schedule, multiply the learning rate by 0.1 every 30 epochs, and train a total of 90 epochs, with a batch size of 256.

Hyperparameters of CIFAR10 For each optimizer, we do a grid search over the learning rate and weight decay for the best hyperparameters. For ADAM and LAPROP, we set $\beta = 0.999$. For MADAM and LAMADAM, we set $\underline{\beta} = 0.5$ and $\bar{\beta} = 0.999$ in all cases. Except for SGD, we tried learning rates from $\{5e-4, 1e-3, 2e-3, 3e-3, 4e-3, 6e-3, 8e-3\}$ and weight decay from $\{0.025, 0.05, 0.1, 0.2, 0.4, 0.8, 1\}$. The best learning rate and weight decay for ADAM, LAPROP, MADAM and LAMADAM are $(3e-3, 0.2)$, $(1e-3, 0.4)$, $(8e-3, 0.05)$ and $(6e-3, 0.05)$ respectively. As to SGD, we tried learning rates from $\{3e-2, 5e-2, 1e-1, 2e-1, 3e-1\}$ and weight decays from $\{1e-4, 3e-4, 5e-4, 1e-3, 2e-3\}$, and the best result was achieved with learning rate $2e-1$ and weight decay $3e-4$. These hyperparameters that gave the best results are also the hyperparameters we used for plotting Figure 7.

Hyperparameters for CIFAR100 We use the same grid search configurations as for CIFAR10. The best learning rate and weight decay for ADAM, LAPROP, MADAM and LAMADAM are $(2e-3, 0.4)$, $(5e-4, 1)$, $(4e-3, 0.2)$ and $(3e-3, 0.2)$ respectively. For SGD, the best learning rate and weight decay are $3e-2$ and $2e-3$ respectively.

Hyperparameters for ImageNet Due to the heavy workload and the time limit, we were not able to accomplish 4 runs for each hyperparameter in ImageNet, so we report the best results for each

⁴<https://github.com/kuangliu/pytorch-cifar>

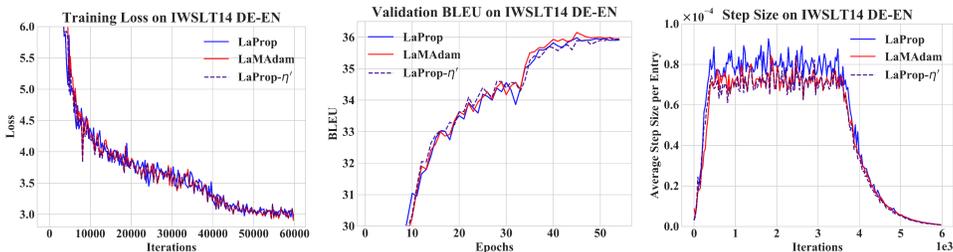


Figure 8: Training loss, validation BLEU and average step size on IWSLT’14 DE-EN, trained with $\eta=5e-4$, $\lambda=1e-2$, $\beta=0.999$ for LAPROP and $\eta=1.5e-3$, $\lambda=1e-2$, $\underline{\beta}=0.5$, $\bar{\beta}=0.999$ for LAMADAM, and $\eta=4.375e-4$, $\lambda=1e-2$, $\beta=0.999$ for LAPROP- η' .

optimizer in Table 5.1, except for the result of ADAM, which was copied from Liu et al. (2020) but uses the same hyperparameters except for the learning rate and weight decay. For LAPROP, MADAM and LAMADAM, we choose learning rates from $\{1e-3, 2e-3, 3e-3, 4e-3, 5e-3, 6e-3, 8e-3\}$ and weight decay from $\{0.003, 0.006, 0.01, 0.012, 0.02, 0.03\}$, and found the best combinations for LAPROP, MADAM and LAMADAM are $(2e-3, 0.03)$, $(5e-3, 0.012)$ and $(6e-3, 0.012)$. For SGD, we choose learning rate from $\{0.05, 0.1, 0.2\}$ and weight decay from $\{5e-5, 7e-5, 1e-4\}$, and found the best combination to be $(0.1, 7e-5)$.

H ADDITIONAL EXPERIMENTAL RESULTS AND DETAILS ON MACHINE TRANSLATION

Additional Experimental Results and analysis In Figure 8, we plot the training loss, validation BLEU and average step size on IWSLT’14 DE-EN. Although the average update size of LAMADAM is smaller even when using 3 times higher learning rate than ADAM, LAMADAM shows slightly better convergence on the training set and better validation BLEU. This may be explained by the heavy-tailed distribution of the gradient in the process of training transformers from scratch Zhang et al. (2019b). Smaller step sizes mitigate the effect of extreme gradient values on the model’s performance. It is worth mentioning that LAPROP diverges using the large learning rate $1.5e-3$. Further, we find LAPROP is unable to produce the same result as LAMADAM even when their update sizes are similar. LAPROP produces a similar step size curve as LAMADAM with learning rate $4.375e-4$, but the performance is weaker than LAMADAM. LAMADAM uses the maximum variation rule to select the adaptive learning rate for each dimension, creating benefit that is not achievable by simply scaling the base learning rate η .

Hyperparameters for IWSLT’14 The transformer we use has 512-dimensional word embeddings and 6 Transformer blocks with 4 attention heads and 1024 FFN dimensions for the encoder/decoder, which is referred to as `transformer_iwslt_de_en` in fairseq. We do a grid search for the learning rate and weight decay for both optimizers. We tried η from $\{2.5e-4, 5e-4, 1e-3, 1.5e-3, 2e-3\}$, and weight decay from $\{0.0001, 0.001, 0.01, 0.1\}$. The best combinations for LAPROP and LAMADAM are $(5e-4, 0.01)$ and $(1.5e-3, 0.01)$. To demonstrate the full potential of adaptive methods under constant learning rates, we use the tri-stage learning rate schedule (Park et al., 2019), linearly increase the learning rate from 0.01η to the full learning rate η in 4k iterations, hold it at η for 32k iterations, and exponentially decay it to 0.01η in 24k iterations. For LAPROP, we tried β from $\{0.98, 0.99, 0.997, 0.999\}$. We found 0.999 to work the best and used it for all the grid search experiments. For LAMADAM, we set $\beta = 0.5$, $\bar{\beta} = 0.999$. For other hyperparameters, we use the default setting in the fairseq example, which sets dropout probability to 0.3, uses label smoothed cross entropy loss with a smoothing coefficient 0.1, and shares the input and output token embedding parameters.

Hyperparameters for WMT’16 The Transformer we use has 1024-dimensional word embeddings, 6 transformer blocks with 16 attention heads and 4096 FFN dimensions for the encoder/decoder, and is referred to as `transformer_vaswani_wmt_en_de_big` in fairseq. The default implementation from fairseq did not use weight decay, so we also ignore weight decay in all experiments. The learning rate schedule takes the first 4k steps to linearly increase the learning

rate to its maximum value. For LAPROP, we found $\beta = 0.98$ to give the best results, and we set $\underline{\beta} = 0.95, \bar{\beta} = 0.98$ in all experiments. This takes around 8 hours on 16 V100 GPUs each run. For grid search, we tried η from $\{5e-4, 1e-3, 1.5e-3, 2e-3\}$, and found $1e-3$ and $1.5e-3$ to work the best for LAPROP and LAMADAM respectively. Other hyperparameters are the defaults of the corresponding fairseq example, which uses a dropout probability of 0.3, the label smoothed cross entropy loss with a smoothing coefficient 0.1, and shares all embedding parameters.

I ADDITIONAL DETAILS OF EXPERIMENTS ON THE GLUE BENCHMARK

The GLUE benchmark is a collection of 9 natural language understanding tasks, namely Corpus of Linguistic Acceptability (CoLA; Warstadt et al. (2018)), Stanford Sentiment Treebank (SST; Socher et al. (2013)), Microsoft Research Paraphrase Corpus (MRPC; Dolan & Brockett (2005)), Semantic Textual Similarity Benchmark (STS; Agirre et al. (2007)), Quora Question Pairs (QQP; Iyer et al. (2017)), Multi-Genre NLI (MNLI; Williams et al. (2018)), Question NLI (QNLI; Rajpurkar et al. (2016)), Recognizing Textual Entailment (RTE; Dagan et al. (2006); Haim et al. (2006); Giampiccolo et al. (2007); Bentivogli et al. (2009)) and Winograd NLI (WNLI; Levesque et al. (2011)).

It is reported in Liu et al. (2019) that ADAM is sensitive to the choice of ϵ on GLUE. Following their settings, we set $\epsilon = 1e - 6$ for both ADAM and MADAM. For LAPROP and LAMADAM, however, we always set $\epsilon = 1e - 15$, like all other experiments in this paper, which is consistent with the observation in Ziyin et al. (2020) that LAPROP is robust to the choice of ϵ . We set $\beta = 0.98$ for ADAM and LAPROP, and $\underline{\beta} = 0.5, \bar{\beta} = 0.98$ for LAPROP and LAMADAM. All other hyperparameters are set to the same as the example in fairseq.⁵ For each task, we do a grid search over the learning rate and weight decay, which are chosen from $\{5e-6, 1e-5, 2e-5, 4e-5, 5e-5, 6e-5\}$ and $\{0.025, 0.05, 0.1, 0.2\}$ respectively. We list the best combinations for ADAM, MADAM, LAPROP and LAMADAM on each task as below:

MNLI: (1e-5, 0.1), (1e-5, 0.1), (4e-5, 0.025), (4e-5, 0.025).

QQP: (1e-5, 0.1), (1e-5, 0.1), (4e-5, 0.025), (4e-5, 0.025).

QNLI: (1e-5, 0.1), (1e-5, 0.1), (4e-5, 0.05), (4e-5, 0.05).

SST-2: (1e-5, 0.1), (1e-5, 0.1), (4e-5, 0.1), (4e-5, 0.1).

RTE: (2e-5, 0.1), (2e-5, 0.1), (6e-5, 0.1), (6e-5, 0.1).

MRPC: (1e-5, 0.1), (1e-5, 0.1), (6e-5, 0.1), (6e-5, 0.1).

STS-B: (2e-5, 0.1), (2e-5, 0.1), (4e-5, 0.5), (4e-5, 0.5).

CoLA: (2e-5, 0.1), (2e-5, 0.1), (6e-5, 0.5), (6e-5, 0.5).

⁵https://github.com/pytorch/fairseq/blob/master/examples/roberta/README_glue.md