
MARS is generalised Nesterov with longer lookahead

Anonymous Author(s)

Affiliation

Address

email

1 MARS (Yuan et al., 2024) is clearly great. But the justification given in the paper makes no sense.
2 Specifically, the justification is given in terms of reducing the minibatch variance by computing the
3 gradient for the same minibatch at different values of the parameters. This would indeed reduce
4 variance. However, this doubles the compute requirements, because you now need to compute the
5 gradient for each datapoint at two setting of the parameters. Instead, they just use the current and
6 previous gradients, evaluated on *different* minibatches. There is no sense in which this can reduce
7 minibatch gradients. Instead, here we argue that MARS should be understood as a generalisation of
8 Nesterov, where you lookahead further than usual.

9 Nesterov is (Sutskever et al., 2013, from)

$$v_{t+1} = \beta_1 v_t - \varepsilon f(\theta_t + \mu v_t) \quad (1a)$$

$$\theta_{t+1} = \theta_t + v_{t+1} \quad (1b)$$

10 The usual interpretation of Nesterov is that you compute the gradient, $f(\theta_t + \mu v_t)$, at a “lookahead”
11 location, $\theta_t + \mu v_t$. Note that in standard Nesterov, the “lookahead” is one momentum step (i.e.
12 $\beta_1 = \mu$). We have generalised Nesterov slightly, by allowing an arbitrary lookahead of size μ .

13 Now, to reframe this to look more like MARS, we start by noting that the gradient steps are performed
14 at $\theta_t + \mu v_t$. We therefore write the v update as,

$$v_{t+1} = \beta_1 v_t - \varepsilon f(\theta'_t) \quad (2)$$

15 which is equivalent to the Nesterov momentum update (Eq. 1a) if we define,

$$\theta'_t = \theta_t + \mu v_t. \quad (3)$$

16 To write the updates for θ' , we take Eq. (3) for timestep $t + 1$ and substitute Eq. (1b),

$$\theta'_{t+1} = \theta_{t+1} + \mu v_{t+1} = \theta_t + v_{t+1} + \mu v_{t+1} \quad (4)$$

17 Adding and subtracting μv_t ,

$$\theta'_{t+1} = \theta_t + \mu v_t - \mu v_t + v_{t+1} + \mu v_{t+1} \quad (5)$$

18 Noticing that $\theta'_t = \theta_t + \mu v_t$ (Eq. 3),

$$\theta'_{t+1} = \theta'_t + v_{t+1} + \mu(v_{t+1} - v_t) \quad (6)$$

19 Thus, the overall updates become,

$$v_{t+1} = \beta_1 v_t - \varepsilon g_t \quad (7a)$$

$$\theta'_{t+1} = \theta'_t + v_{t+1} + \mu(v_{t+1} - v_t) \quad (7b)$$

20 where $g_t = f(\theta'_t)$ is the gradient evaluated at θ'_t .

21 Here, the $v_{t+1} - v_t$ correction term arises in the parameter update, not in the parameter update, as in
22 MARS. We can push the correction into the momentum by writing the parameter update as,

$$\theta'_{t+1} = \theta'_t + v'_{t+1} \quad (8)$$

23 which is equivalent to the Nesterov parameter update (Eq. 7b) if we define,

$$v'_{t+1} = v_{t+1} + \mu(v_{t+1} - v_t). \quad (9)$$

24 To write the updates for v' , we substitute the update for v (Eq. (7a), into Eq. 9,

$$v'_{t+1} = (\beta_1 v_t - \varepsilon g_t) + \mu((\beta_1 v_t - \varepsilon g_t) - (\beta_1 v_{t-1} - \varepsilon g_{t-1})) \quad (10)$$

25 Rearranging,

$$v'_{t+1} = \beta_1(v_t + \mu(v_t - v_{t-1})) - \varepsilon(g_t + \mu(g_t - g_{t-1})) \quad (11)$$

26 Identifying $v'_t = v_t + \mu(v_t - v_{t-1})$ (Eq. 9),

$$v'_{t+1} = \beta_1 v'_t - \varepsilon(g_t + \mu(g_t - g_{t-1})). \quad (12)$$

27 Overall, the updates become,

$$\theta'_{t+1} = \theta'_t + v'_{t+1} \quad (13a)$$

$$v'_{t+1} = \beta_1 v'_t - \varepsilon(g_t + \mu(g_t - g_{t-1})). \quad (13b)$$

28 This is starting to resemble MARS! Critically, in MARS, there is a learning rate that applies to the
29 gradient in the parameter update. To get something like this in our setup, we define v'' as a scaled
30 version of v' ,

$$\theta'_{t+1} = \theta'_t + \eta v''_{t+1} \quad (14)$$

31 This is equivalent to Eq. (13b) if we set,

$$v''_{t+1} = \frac{1}{\eta} v'_{t+1} \quad (15)$$

32 To get updates for v'' , we substitute Eq. 13b,

$$v''_{t+1} = \frac{1}{\eta} \beta_1 v'_t - \frac{\varepsilon}{\eta} (g_t + \mu(g_t - g_{t-1})). \quad (16)$$

33 Recognising that $v'_t/\eta = v''_t$ (Eq. 15), we get,

$$v''_{t+1} = \beta_1 v''_t - \frac{\varepsilon}{\eta} (g_t + \mu(g_t - g_{t-1})). \quad (17)$$

34 Overall, this gives updates,

$$\theta'_{t+1} = \theta'_t + \eta v''_{t+1} \quad (18a)$$

$$v''_{t+1} = \beta_1 v''_t - \frac{\varepsilon}{\eta} (g_t + \mu(g_t - g_{t-1})). \quad (18b)$$

35 which has exactly the same form as MARS, except for the constants. Specifically, the MARS updates
36 are the same for θ' , but have different constants for v'' ,

$$v''_{t+1} = \beta_1 v''_t - (1 - \beta_1)g_t + \beta_1 \mu(g_t - g_{t-1}). \quad (19)$$

37 Thus, MARS can be connected to generalised Nesterov by taking the constants in Eq. 18b and Eq. 19
38 to be equal. Specifically, taking the constants for the g_t terms to be equal,

$$\frac{\varepsilon}{\eta} = 1 - \beta \quad (20)$$

implies that to make generalised Nesterov equivalent to MARS, we need to set the generalised
39 Nesterov parameter, ε , to,

$$\varepsilon = \eta(1 - \beta_1) \quad (21)$$

40 And, taking the constants for the $\mu(g_t - g_{t-1})$ terms to be equal,

$$\frac{\varepsilon}{\eta} \mu = \beta_1 \quad (22)$$

implies that to make generalized Nesterov equivalent to MARS, we need to set the gearlised
41 Nesterov parameter, μ , to,

$$\mu = \frac{\beta_1}{\frac{\varepsilon}{\eta}} = \frac{\beta_1}{1 - \beta_1} \quad (23)$$

42 What does this mean for the interpretation of MARS? It means that MARS is generalised Nesterov,
43 where we lookahead a long way: specifically, the look-ahead is set by μ . In standard Nesterov, we
44 would have $\mu = \beta_1$, and remember that β_1 is typically close to 1, e.g. $\mu = \beta_1 = 0.9$. However, in
45 MARS, with $\beta_1 = 0.9$, the lookahead is given by,

$$\mu(\beta_1 = 0.9) = \frac{0.9}{1 - 0.1} = \frac{0.9}{1 - 0.1} = \frac{0.9}{0.1} = 9. \quad (24)$$

46 Implying that if $\beta_1 = 0.9$, MARS is generalised Nesterov which looks ahead ten times further than
47 standard Nesterov.

48 **References**

- 49 Sutskever, I., Martens, J., Dahl, G., and Hinton, G. On the importance of initialization and momentum
50 in deep learning. In *International conference on machine learning*, pp. 1139–1147. PMLR, 2013.
- 51 Yuan, H., Liu, Y., Wu, S., Zhou, X., and Gu, Q. Mars: Unleashing the power of variance reduction
52 for training large models. *arXiv:2411.10438*, 2024.