

Gradient-based Hyperparameter Optimization through Reversible Learning

Dougal Maclaurin, David Duvenaud, Ryan P. Adams

CSC 2541 presentation by
Haoping Xu, Zhihuan Yu, Jingcheng Niu

Overview

- The bilevel hyperparameter optimization of SGD with momentum
- How to calculate the gradient of the inner SGD
- Implementation detail: precision
- Experiments
 - Dataset distillation
 - Learning rate scheduling

Hyperparameter Optimization of SGD with momentum

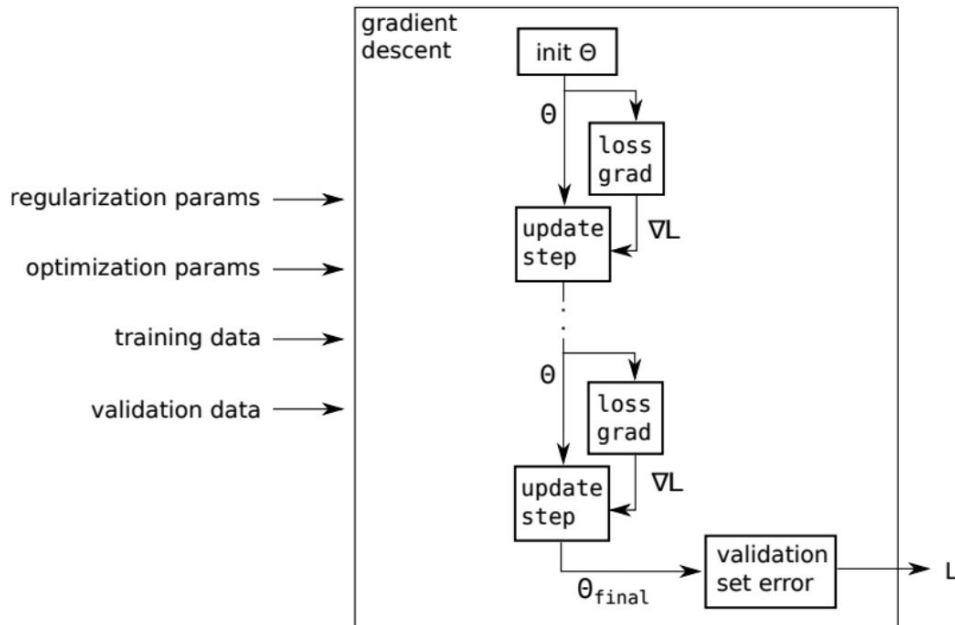
Inner objective:

SGD with momentum

Algorithm 1 Stochastic gradient descent with momentum

- 1: **input:** initial \mathbf{w}_1 , decays γ , learning rates α , loss function $L(\mathbf{w}, \boldsymbol{\theta}, t)$
 - 2: initialize $\mathbf{v}_1 = \mathbf{0}$
 - 3: **for** $t = 1$ **to** T **do**
 - 4: $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$ ▷ evaluate gradient
 - 5: $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$ ▷ update velocity
 - 6: $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$ ▷ update position
 - 7: **end for**
 - 8: **output** trained parameters \mathbf{w}_T
-

Unrolling



Reverse-mode differentiation (RMD) of SGD

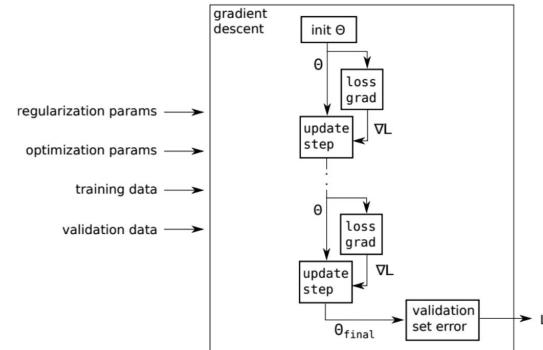
Algorithm 2 Reverse-mode differentiation of SGD

- 1: **input:** $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$, train loss $L(\mathbf{w}, \boldsymbol{\theta}, t)$, loss $f(\mathbf{w})$
- 2: initialize $d\mathbf{v} = \mathbf{0}, d\boldsymbol{\theta} = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$
- 3: initialize $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$
- 4: **for** $t = T$ **counting down to 1 do**
- 5: $d\alpha_t = d\mathbf{w}^\top \mathbf{v}_t$
- 6: $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$
- 7: $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
- 8: $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$
- 9: $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$
- 10: $d\gamma_t = d\mathbf{v}^\top (\mathbf{v}_t + \mathbf{g}_t)$
- 11: $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
- 12: $d\boldsymbol{\theta} = d\boldsymbol{\theta} - (1 - \gamma_t) d\mathbf{v} \nabla_{\boldsymbol{\theta}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
- 13: $d\mathbf{v} = \gamma_t d\mathbf{v}$
- 14: **end for**
- 15: **output** gradient of $f(\mathbf{w}_T)$ w.r.t $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$ and $\boldsymbol{\theta}$

} exactly reverse
gradient descent
operations

Algorithm 1 Stochastic gradient descent with momentum

- 1: **input:** initial \mathbf{w}_1 , decays γ , learning rates α , loss function $L(\mathbf{w}, \boldsymbol{\theta}, t)$
- 2: initialize $\mathbf{v}_1 = \mathbf{0}$
- 3: **for** $t = 1$ **to** T **do**
- 4: $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$ ▷ evaluate gradient
- 5: $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$ ▷ update velocity
- 6: $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$ ▷ update position
- 7: **end for**
- 8: **output** trained parameters \mathbf{w}_T



Reverse-mode differentiation (RMD) of SGD

Algorithm 2 Reverse-mode differentiation of SGD

- 1: **input:** $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$, train loss $L(\mathbf{w}, \boldsymbol{\theta}, t)$, loss $f(\mathbf{w})$
 - 2: initialize $d\mathbf{v} = \mathbf{0}, d\boldsymbol{\theta} = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$
 - 3: initialize $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$
 - 4: **for** $t = T$ **counting down to 1 do**
 - 5: $d\alpha_t = d\mathbf{w}^\top \mathbf{v}_t$
 - 6: $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$
 - 7: $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
 - 8: $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$
 - 9: $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$
 - 10: $d\gamma_t = d\mathbf{v}^\top (\mathbf{v}_t + \mathbf{g}_t)$
 - 11: $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
 - 12: $d\boldsymbol{\theta} = d\boldsymbol{\theta} - (1 - \gamma_t) d\mathbf{v} \nabla_{\boldsymbol{\theta}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
 - 13: $d\mathbf{v} = \gamma_t d\mathbf{v}$
 - 14: **end for**
 - 15: **output** gradient of $f(\mathbf{w}_T)$ w.r.t $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$ and $\boldsymbol{\theta}$
-

Quite a few of gradients!

Reverse-mode differentiation (RMD) of SGD

Algorithm 2 Reverse-mode differentiation of SGD

- 1: **input:** $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$, train loss $L(\mathbf{w}, \boldsymbol{\theta}, t)$, loss $f(\mathbf{w})$
 - 2: initialize $d\mathbf{v} = \mathbf{0}, d\boldsymbol{\theta} = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$
 - 3: initialize $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$
 - 4: **for** $t = T$ **counting down to 1** **do**
 - 5: $d\alpha_t = d\mathbf{w}^\top \mathbf{v}_t$
 - 6: $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$
 - 7: $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
 - 8: $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$
 - 9: $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$
 - 10: $d\gamma_t = d\mathbf{v}^\top (\mathbf{v}_t + \mathbf{g}_t)$
 - 11: $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
 - 12: $d\boldsymbol{\theta} = d\boldsymbol{\theta} - (1 - \gamma_t) d\mathbf{v} \nabla_{\boldsymbol{\theta}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$
 - 13: $d\mathbf{v} = \gamma_t d\mathbf{v}$
 - 14: **end for**
 - 15: **output** gradient of $f(\mathbf{w}_T)$ w.r.t $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$ and $\boldsymbol{\theta}$
-

Autograd comes in handy!

loss_grad = grad(loss)

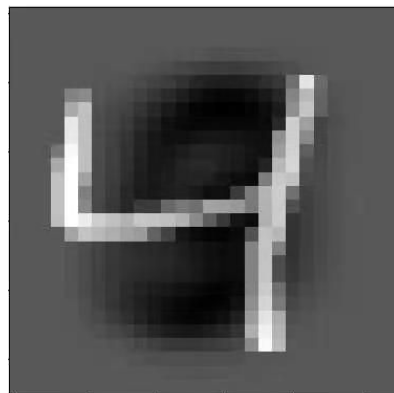
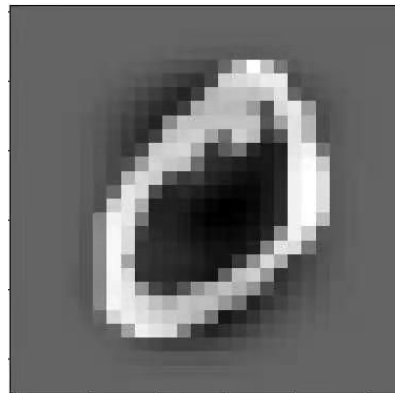
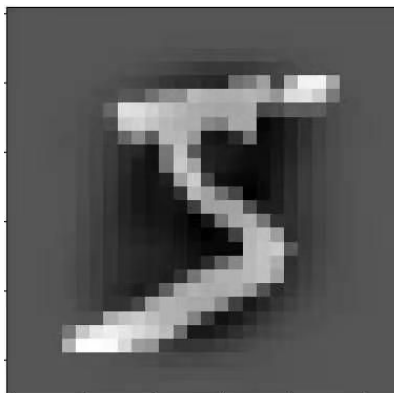
loss_hvp = grad(loss_grad)

meta_hvp = grad(loss_grad, argnums=1)

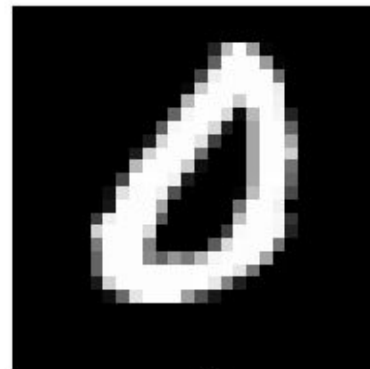
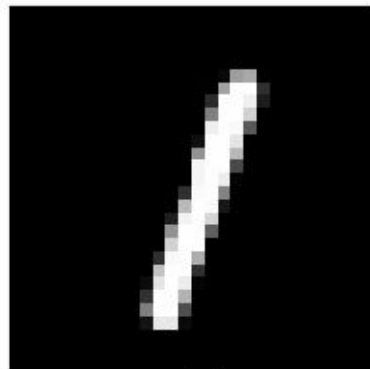
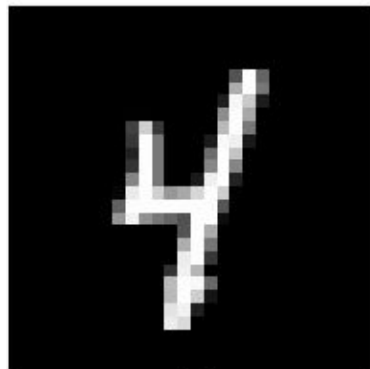
} exactly reverse
} gradient descent
} operations

Side note: MNIST data

The paper's
MNIST



Original
MNIST



Numerical Precision

In practice, Algorithm 2 fails due to finite numerical precision.

$$5: \quad \mathbf{v}_t = \gamma \mathbf{v}_{t-1} + (1 - \gamma) \mathbf{g}_{t-1}$$

$$8: \quad \mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$$

The reverse process requires repeated multiplication by $1/\gamma$, causing errors accumulate exponentially!

Exactly Reversible Multiplication

So instead, we represent v and w each with a 64-bit integer and a remainder buffer.

Using this method we have can achieve float64 precision and our computations is reversible.

Algorithm 3 Exactly reversible multiplication by a ratio

- 1: **Input:** Information buffer i , value c , ratio n/d
 - 2: $i = i \times d$ ▷ make room for new digit
 - 3: $i = i + (c \bmod d)$ ▷ store digit lost by division
 - 4: $c = c \div d$ ▷ divide by denominator
 - 5: $c = c \times n$ ▷ multiply by numerator
 - 6: $c = c + (i \bmod n)$ ▷ add digit from buffer
 - 7: $i = i \div n$ ▷ shorten information buffer
 - 8: **return** updated buffer i , updated value c
-

Experiments

- Dataset distillation
- Learning rate scheduling

The logo for Google Colab, featuring the word "colab" in a bold, lowercase, sans-serif font. The letters "c" and "o" are yellow with a white outline, while the letters "l", "a", and "b" are solid orange.