

Towards Understanding Why LookAhead Generalizes Better Than SGD and Beyond

Pan Zhou

Joint work with Hanshu Yan, Xiaotong Yuan, Jiashi Feng, and Shuicheng Yan

SEA AI Lab

zhoupan@sea.com

Dec 06, 2021

Background: LookAhead

Algorithm 1: SGD

Input : Objective $F_S(\theta)$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.

```
for  $t = 1, 2, \dots, T$  do
   $\mathbf{v}_0^{(t)} = \theta_{t-1}$ ;
  for  $\tau = 1, 2, \dots, k$  do
     $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$ 
  end
   $\theta_t = \mathbf{v}_k^{(t)} = (1 - 1)\theta_{t-1} + 1 * \mathbf{v}_k^{(t)}$  ( $\alpha = 1$ )
end
```

outer-loop optimization

Output : $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

Algorithm 2: Lookahead

Input : Objective $F_S(\theta)$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.

```
for  $t = 1, 2, \dots, T$  do
   $\mathbf{v}_0^{(t)} = \theta_{t-1}$ ;
  for  $\tau = 1, 2, \dots, k$  do
     $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$ 
  end
   $\theta_t = (1 - \alpha)\theta_{t-1} + \alpha \mathbf{v}_k^{(t)}$ .
end
```

outer-loop optimization

Output : $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

[1] Lookahead Optimizer: k steps forward, 1 step back, NeurIPS'19

Background: LookAhead

Algorithm 1: SGD

Input : Objective $F_S(\theta)$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.

```
for  $t = 1, 2, \dots, T$  do
   $\mathbf{v}_0^{(t)} = \theta_{t-1}$ ;
  Inner-loop optimization
  for  $\tau = 1, 2, \dots, k$  do
     $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$ 
  end
   $\theta_t = \mathbf{v}_k^{(t)} = (1 - 1)\theta_{t-1} + 1 * \mathbf{v}_k^{(t)}$  ( $\alpha = 1$ )
```

end
Output : $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

Algorithm 2: Lookahead

Input : Objective $F_S(\theta)$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.

```
for  $t = 1, 2, \dots, T$  do
   $\mathbf{v}_0^{(t)} = \theta_{t-1}$ ;
  Inner-loop optimization
  for  $\tau = 1, 2, \dots, k$  do
     $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$ 
  end
   $\theta_t = (1 - \alpha)\theta_{t-1} + \alpha \mathbf{v}_k^{(t)}$ .
```

end
Output : $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

Key steps in SGD & LookAhead (LA):

inner-loop optimization: K steps forward in SGD & LA

outer-loop optimization: 1 step back in LA, while no step back in SGD

[1] Lookahead Optimizer: k steps forward, 1 step back, NeurIPS'19

Background: LookAhead

Algorithm 1: SGD

Input : Objective $F_S(\theta)$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.

```
for  $t = 1, 2, \dots, T$  do
   $\mathbf{v}_0^{(t)} = \theta_{t-1}$ ;
  for  $\tau = 1, 2, \dots, k$  do
     $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$ 
  end
   $\theta_t = \mathbf{v}_k^{(t)} = (1 - 1)\theta_{t-1} + 1 * \mathbf{v}_k^{(t)}$  ( $\alpha = 1$ )
end
```

outer-loop optimization

Output : $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

Algorithm 2: Lookahead

Input : Objective $F_S(\theta)$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.

```
for  $t = 1, 2, \dots, T$  do
   $\mathbf{v}_0^{(t)} = \theta_{t-1}$ ;
  for  $\tau = 1, 2, \dots, k$  do
     $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$ 
  end
   $\theta_t = (1 - \alpha)\theta_{t-1} + \alpha \mathbf{v}_k^{(t)}$ .
end
```

outer-loop optimization

Output : $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

Key steps in SGD & LookAhead (LA):

inner-loop optimization: K steps forward in SGD & LA

outer-loop optimization: 1 step back in LA, while no step back in SGD

[1] Lookahead Optimizer: k steps forward, 1 step back, NeurIPS'19

Observation: LookAhead Generalizes Better Than SGD

Important observations: LookAhead (LA) [1] enjoys better test performance than SGD

OPTIMIZER	CIFAR-10	CIFAR-100
SGD	95.23 \pm .19	78.24 \pm .18
POLYAK	95.26 \pm .04	77.99 \pm .42
ADAM	94.84 \pm .16	76.88 \pm .39
LOOKAHEAD	95.27 \pm .06	78.34 \pm .05

ResNet 18 [1]

OPTIMIZER	TRAIN	VAL.	TEST
SGD	43.62	66.0	63.90
LA(SGD)	35.02	65.10	63.04
ADAM	33.54	61.64	59.33
LA(ADAM)	31.92	60.28	57.72
POLYAK	-	61.18	58.79

LSTM [1]

[1] Lookahead Optimizer: k steps forward, 1 step back, NeurIPS'19

Observation: LookAhead Generalizes Better Than SGD

Important observations: LookAhead (LA) [1] enjoys better test performance than SGD

OPTIMIZER	CIFAR-10	CIFAR-100
SGD	95.23 \pm .19	78.24 \pm .18
POLYAK	95.26 \pm .04	77.99 \pm .42
ADAM	94.84 \pm .16	76.88 \pm .39
LOOKAHEAD	95.27 \pm .06	78.34 \pm .05

ResNet 18 [1]

OPTIMIZER	TRAIN	VAL.	TEST
SGD	43.62	66.0	63.90
LA(SGD)	35.02	65.10	63.04
ADAM	33.54	61.64	59.33
LA(ADAM)	31.92	60.28	57.72
POLYAK	-	61.18	58.79

LSTM [1]

Problems:

1. Why LA enjoys better test performance than SGD?
2. How to further improve LA?

[1] Lookahead Optimizer: k steps forward, 1 step back, NeurIPS'19

Tools for Test Performance Analysis

- **optimal solution** to empirical risk:

$$\boldsymbol{\theta}_S^* \in \operatorname{argmin}_{\boldsymbol{\theta}} F_S(\boldsymbol{\theta}) \triangleq \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \boldsymbol{\theta}), \mathbf{y}_i), \quad (1)$$

- **approximate solution** to empirical risk:

$$\boldsymbol{\theta}_{A,S} \approx \operatorname{argmin}_{\boldsymbol{\theta}} F_S(\boldsymbol{\theta}) \triangleq \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \boldsymbol{\theta}), \mathbf{y}_i), \quad (2)$$

Tools for Test Performance Analysis

- **optimal solution** to empirical risk:

$$\boldsymbol{\theta}_S^* \in \operatorname{argmin}_{\boldsymbol{\theta}} F_S(\boldsymbol{\theta}) \triangleq \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \boldsymbol{\theta}), \mathbf{y}_i), \quad (1)$$

- **approximate solution** to empirical risk:

$$\boldsymbol{\theta}_{\mathcal{A},S} \approx \operatorname{argmin}_{\boldsymbol{\theta}} F_S(\boldsymbol{\theta}) \triangleq \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \boldsymbol{\theta}), \mathbf{y}_i), \quad (2)$$

- **Excess risk error** to measure **test performance**:

excess risk error

$$\varepsilon_{\text{exc}} = \underbrace{\mathbb{E}_{\mathcal{A},S}[F(\boldsymbol{\theta}_{\mathcal{A},S})]}_{\text{test error}} - \underbrace{\mathbb{E}_{\mathcal{A},S}[F_S(\boldsymbol{\theta}_S^*)]}_{\text{best training error}} = \underbrace{\mathbb{E}_{\mathcal{A},S}[F(\boldsymbol{\theta}_{\mathcal{A},S}) - F_S(\boldsymbol{\theta}_{\mathcal{A},S})]}_{\text{generalization error}} + \underbrace{\mathbb{E}_{\mathcal{A},S}[F_S(\boldsymbol{\theta}_{\mathcal{A},S}) - F_S(\boldsymbol{\theta}_S^*)]}_{\text{optimization error}} \quad (3)$$

where $F(\boldsymbol{\theta}) \triangleq \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[\ell(f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y})]$ is population risk.

Superiority of LA: Smaller Optimization & Generalization Errors

Excess risk error to measure test performance:

$$\varepsilon_{\text{exc}} = \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}})]}_{\text{test error}} - \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}^*)]}_{\text{best training error}} = \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}}) - F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}})]}_{\text{generalization error}} + \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}}) - F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}^*)]}_{\text{optimization error}}$$

Theorem 1 (informal) Under proper assumptions, by setting conventional learning rate $\eta = 1/\sqrt{kT}$, on **convex problem** we have

$$\text{optimization error} \leq O\left(\frac{1}{\alpha\sqrt{kT}}\right) \qquad \text{generalization error} \leq O\left(\frac{\alpha\sqrt{kT}}{n}\right) \qquad (5)$$

where kT is total training iteration number, n is training sample number.

Superiority of LA: Smaller Optimization & Generalization Errors

Excess risk error to measure test performance:

$$\varepsilon_{\text{exc}} = \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}})]}_{\text{test error}} - \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}^*)]}_{\text{best training error}} = \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}}) - F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}})]}_{\text{generalization error}} + \underbrace{\mathbb{E}_{\mathcal{A},\mathcal{S}}[F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{A},\mathcal{S}}) - F_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}^*)]}_{\text{optimization error}}$$

Theorem 1 (informal) Under proper assumptions, by setting conventional learning rate $\eta = 1/\sqrt{kT}$, on **convex problem** we have

$$\text{optimization error} \leq O\left(\frac{1}{\alpha\sqrt{kT}}\right) \qquad \text{generalization error} \leq O\left(\frac{\alpha\sqrt{kT}}{n}\right) \qquad (5)$$

where kT is total training iteration number, n is training sample number.

Since (1) optimum of α is $\alpha = O(1 \wedge \sqrt{n/kT})$ and (2) SGD = LA with $\alpha = 1$

Lookahead enjoys smaller excess risk error (test error) than SGD

Superiority of LA: Smaller Optimization & Generalization Errors

Theorem 2 (informal). λ -**strongly-convex problem** with proper assumptions:

$$\text{optimization error} \leq \begin{cases} O\left(\frac{1}{T^{2\alpha}} + \frac{1}{\lambda^2(kT)^{2\alpha(1-2\alpha)}}\right), & 0 < \alpha < \frac{1}{2}, \\ O\left(\frac{1}{T} + \frac{\log(Tk)}{\lambda^2 k T}\right), & \alpha = \frac{1}{2}, \\ O\left(\frac{1}{T^{2\alpha}} + \frac{1}{\lambda^2(2\alpha-1)kT}\right), & \frac{1}{2} < \alpha \leq 1. \end{cases} \quad \text{generalization error} \leq O\left(\frac{1}{n\lambda} \frac{(Tk+1)^\alpha - 1}{((T+1)k+2)^\alpha}\right) \quad (6)$$

When problem is large-scale and iteration number T is not large,

$$\frac{\ln T}{T^\alpha} > O\left(\frac{1}{n\lambda} \frac{\ln(Tk)}{k^\alpha}\right) \quad \text{or} \quad \frac{1}{\lambda(\alpha-1)^2 T k} > O\left(\frac{1}{n} \frac{\ln(Tk)}{T^\alpha k^\alpha}\right) \quad (7)$$

the optimum α is not 1.

Superiority of LA: Smaller Optimization & Generalization Errors

Theorem 2 (informal). λ -**strongly-convex problem** with proper assumptions:

$$\text{optimization error} \leq \begin{cases} O\left(\frac{1}{T^{2\alpha}} + \frac{1}{\lambda^2(kT)^{2\alpha(1-2\alpha)}}\right), & 0 < \alpha < \frac{1}{2}, \\ O\left(\frac{1}{T} + \frac{\log(Tk)}{\lambda^2 k T}\right), & \alpha = \frac{1}{2}, \\ O\left(\frac{1}{T^{2\alpha}} + \frac{1}{\lambda^2(2\alpha-1)kT}\right), & \frac{1}{2} < \alpha \leq 1. \end{cases} \quad \text{generalization error} \leq O\left(\frac{1}{n\lambda} \frac{(Tk+1)^\alpha - 1}{((T+1)k+2)^\alpha}\right) \quad (6)$$

When problem is large-scale and iteration number T is not large,

$$\frac{\ln T}{T^\alpha} > O\left(\frac{1}{n\lambda} \frac{\ln(Tk)}{k^\alpha}\right) \quad \text{or} \quad \frac{1}{\lambda(\alpha-1)^2 T k} > O\left(\frac{1}{n} \frac{\ln(Tk)}{T^\alpha k^\alpha}\right) \quad (7)$$

the optimum α is not 1.

This also explains **why Lookahead enjoys smaller excess risk error (test error) than SGD**

Superiority of LA: Smaller Optimization & Generalization Errors

Theorem 3 (informal). On nonconvex problem with PL condition

$$2\mu(f(w) - f(w^*)) \leq \|\nabla f(w)\|^2 \quad (\forall w, w^* \in \operatorname{argmin}_w f(w))$$

then we have

$$\text{optimization error} \leq O\left(\frac{1}{(Tk + 1)^{2\alpha}} + \frac{\alpha(\alpha + 2(1 - \alpha)(k - 1))}{\mu^2(Tk + 1)^{2\alpha - 1}}\right)$$

$$\text{generalization error} \leq O\left(\frac{1}{n - 1} \alpha^{\frac{1}{1 + \gamma}} (Tk)^{\frac{\gamma}{\gamma + 1}}\right) \quad \left(\gamma = \left(1 - \frac{1}{n}\right) \frac{\alpha L}{\mu}\right)$$

Superiority of LA: Smaller Optimization & Generalization Errors

Theorem 3 (informal). On nonconvex problem with PL condition

$$2\mu(f(w) - f(w^*)) \leq \|\nabla f(w)\|^2 \quad (\forall w, w^* \in \operatorname{argmin}_w f(w))$$

then we have

$$\begin{aligned} \text{optimization error} &\leq O\left(\frac{1}{(Tk + 1)^{2\alpha}} + \frac{\alpha(\alpha + 2(1 - \alpha)(k - 1))}{\mu^2(Tk + 1)^{2\alpha - 1}}\right) \\ \text{generalization error} &\leq O\left(\frac{1}{n - 1} \alpha^{\frac{1}{1 + \gamma}} (Tk)^{\frac{\gamma}{\gamma + 1}}\right) \quad \left(\gamma = \left(1 - \frac{1}{n}\right) \frac{\alpha L}{\mu}\right) \end{aligned}$$

By properly choosing α ,

Lookahead can also enjoys smaller excess risk error (test error) than SGD

An Improved LookAhead: Stagewise Locally-regularized Lookahead

Algorithm 3: Stagewise Locally-Regularized LookAhead (SLRLA)

Input : Loss $F_S(\boldsymbol{\theta})$, constant $\{\beta_q\}_{q=1}^Q$
for $q = 1, 2, \dots, Q$ **do**
 $F_q(\boldsymbol{\theta}) = F_S(\boldsymbol{\theta}) + \frac{\beta_q}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_{q-1}\|^2$;
 $\boldsymbol{\theta}_q = \text{Look-ahead}(F_q(\boldsymbol{\theta}), \eta_q, T_q, \alpha_q, k_q, \boldsymbol{\theta}_{q-1}, \mathcal{A}, \mathcal{S})$. 
end
Output : $\boldsymbol{\theta}_{\mathcal{A}, \mathcal{S}} = \boldsymbol{\theta}_Q$.

Algorithm 2: Lookahead

Input : Objective $F_S(\boldsymbol{\theta})$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.
for $t = 1, 2, \dots, T$ **do**
 $\mathbf{v}_0^{(t)} = \boldsymbol{\theta}_{t-1}$;
 for $\tau = 1, 2, \dots, k$ **do**
 $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\boldsymbol{\theta}), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$
 end
 $\boldsymbol{\theta}_t = (1 - \alpha)\boldsymbol{\theta}_{t-1} + \alpha \mathbf{v}_k^{(t)}$.
end
Output : $\boldsymbol{\theta}_{\mathcal{A}, \mathcal{S}} = \boldsymbol{\theta}_T$

Strategy: divide optimization into **several stages** and use lookahead to solve **locally-regularized loss**

$$F_q(\boldsymbol{\theta}) = F_S(\boldsymbol{\theta}) + \frac{\beta_q}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_{q-1}\|^2$$

An Improved LookAhead: Stagewise Locally-regularized Lookahead

Algorithm 3: Stagewise Locally-Regularized LookAhead (SLRLA)

Input : Loss $F_S(\boldsymbol{\theta})$, constant $\{\beta_q\}_{q=1}^Q$
for $q = 1, 2, \dots, Q$ **do**
 $F_q(\boldsymbol{\theta}) = F_S(\boldsymbol{\theta}) + \frac{\beta_q}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_{q-1}\|^2$;
 $\boldsymbol{\theta}_q = \text{Look-ahead}(F_q(\boldsymbol{\theta}), \eta_q, T_q, \alpha_q, k_q, \boldsymbol{\theta}_{q-1}, \mathcal{A}, \mathcal{S})$. 
end
Output : $\boldsymbol{\theta}_{\mathcal{A}, \mathcal{S}} = \boldsymbol{\theta}_Q$.

Algorithm 2: Lookahead

Input : Objective $F_S(\boldsymbol{\theta})$, dataset \mathcal{S} , inner-loop optimizer \mathcal{A} , inner-loop step number k and learning rate $\{\{\eta_\tau^{(t)}\}\}$, outer-loop learning rate $\alpha \in (0, 1)$.
for $t = 1, 2, \dots, T$ **do**
 $\mathbf{v}_0^{(t)} = \boldsymbol{\theta}_{t-1}$;
 for $\tau = 1, 2, \dots, k$ **do**
 $\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\boldsymbol{\theta}), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$
 end
 $\boldsymbol{\theta}_t = (1 - \alpha)\boldsymbol{\theta}_{t-1} + \alpha \mathbf{v}_k^{(t)}$.
end
Output : $\boldsymbol{\theta}_{\mathcal{A}, \mathcal{S}} = \boldsymbol{\theta}_T$

Strategy: divide optimization into **several stages** and use lookahead to solve **locally-regularized loss**

$$F_q(\boldsymbol{\theta}) = F_S(\boldsymbol{\theta}) + \frac{\beta_q}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_{q-1}\|^2$$

Advantages:

- Local regularization improves loss convexity, e.g. ill-conditioned loss \rightarrow well-conditioned one
- Local regularization helps avoid overfitting

Superiority of SLRLA: Smaller Optimization & Generalization Errors

Theorem 4 (informal). Under proper assumptions, to obtain optimization error

$$\text{optimization error} \leq \epsilon$$

The stochastic gradient complexity (stochastic gradient evaluation number, a.k.a. IFO) is

stochastic gradient complexity	LookAhead (LA)			SLRLA
	$\alpha \in (0, \frac{1}{2})$	$\alpha = \frac{1}{2}$	$\alpha \in (\frac{1}{2}, 1]$	$\alpha \in (0, 1]$
λ -strongly-convex problems	$\mathcal{O}\left(\left(\frac{1}{\epsilon}\right)^{\frac{1}{2\alpha}} + \left(\frac{1}{(1-2\alpha)\lambda^2\epsilon}\right)^{\frac{1}{2\alpha}}\right)$	$\mathcal{O}\left(\frac{\log \frac{1}{\epsilon}}{\lambda^2\epsilon}\right)$	$\mathcal{O}\left(\frac{1}{(2\alpha-1)\lambda^2\epsilon}\right)$	$\mathcal{O}\left(\frac{1}{\lambda\alpha\epsilon}\right)$
nonconvex problems with μ -PL		$\mathcal{O}\left(\left(\frac{1}{\mu^2\epsilon}\right)^{1/\alpha}\right)$		$\mathcal{O}\left(\frac{1}{\mu\alpha\epsilon}\right)$

Superiority of SLRLA: Smaller Optimization & Generalization Errors

Theorem 4 (informal). Under proper assumptions, to obtain optimization error
 optimization error $\leq \epsilon$

The stochastic gradient complexity (stochastic gradient evaluation number, a.k.a. IFO) is

stochastic gradient complexity	LookAhead (LA)			SLRLA
	$\alpha \in (0, \frac{1}{2})$	$\alpha = \frac{1}{2}$	$\alpha \in (\frac{1}{2}, 1]$	$\alpha \in (0, 1]$
λ -strongly-convex problems	$\mathcal{O}\left(\left(\frac{1}{\epsilon}\right)^{\frac{1}{2\alpha}} + \left(\frac{1}{(1-2\alpha)\lambda^2\epsilon}\right)^{\frac{1}{2\alpha}}\right)$	$\mathcal{O}\left(\frac{\log \frac{1}{\epsilon}}{\lambda^2\epsilon}\right)$	$\mathcal{O}\left(\frac{1}{(2\alpha-1)\lambda^2\epsilon}\right)$	$\mathcal{O}\left(\frac{1}{\lambda\alpha\epsilon}\right)$
nonconvex problems with μ -PL		$\mathcal{O}\left(\left(\frac{1}{\mu^2\epsilon}\right)^{1/\alpha}\right)$		$\mathcal{O}\left(\frac{1}{\mu\alpha\epsilon}\right)$

By observing factors α , λ and μ , SLRLA has smaller computational complexity than LA, meaning

SLRLA has smaller optimization error than LA under a given computational budget

Superiority of SLRLA: Smaller Optimization & Generalization Errors

Theorem 4 (informal). Under proper assumptions, to obtain optimization error

$$\text{optimization error} \leq \epsilon$$

the generalization error is

generalization error	LookAhead (LA) $\alpha \in (0, 1]$	SLRLA $\alpha \in (0, 1]$
λ -strongly-convex problems	$\mathcal{O}\left(\frac{1}{n\lambda}\right)$	$\mathcal{O}\left(\frac{1}{n(\beta/\alpha + \lambda)}\right)$
nonconvex problems with μ -PL	$\mathcal{O}\left(\frac{1}{n} (Tk)^{\frac{\gamma}{\gamma+1}}\right)$ ($\gamma = (1 - \frac{1}{n}) \frac{\alpha L}{\mu}$)	$\mathcal{O}\left(\frac{1}{n} / \left(\frac{c}{\alpha} + \mu\right)\right)$ ($c \geq 0$)

Superiority of SLRLA: Smaller Optimization & Generalization Errors

Theorem 4 (informal). Under proper assumptions, to obtain optimization error

$$\text{optimization error} \leq \epsilon$$

the generalization error is

generalization error	LookAhead (LA) $\alpha \in (0, 1]$	SLRLA $\alpha \in (0, 1]$
λ -strongly-convex problems	$\mathcal{O}\left(\frac{1}{n\lambda}\right)$	$\mathcal{O}\left(\frac{1}{n(\beta/\alpha + \lambda)}\right)$
nonconvex problems with μ -PL	$\mathcal{O}\left(\frac{1}{n} (Tk)^{\frac{\gamma}{\gamma+1}}\right)$ ($\gamma = (1 - \frac{1}{n}) \frac{\alpha L}{\mu}$)	$\mathcal{O}\left(\frac{1}{n} / \left(\frac{c}{\alpha} + \mu\right)\right)$ ($c \geq 0$)

By comparison,

SLRLA has smaller generalization error than LA

Experimental Results

Table 3: Classification accuracy (%). \diamond , $*$, \dagger , \ddagger are respectively reported in [1], [15], [49], [50].

optimizer	CIFAR10			CIFAR100			ImageNet
	ResNet18	VGG16	WRN-16-10	ResNet18	VGG16	WRN-16-10	ResNet18
Adam [11]	94.84 \diamond	91.08	93.54	76.88 \diamond	64.07	74.81	66.54 $*$
Adabound [51]	92.56	91.35	91.68	71.43	64.74	71.64	68.13 \dagger
RAdam [15]	93.85	90.84	94.16	74.30	63.99	75.92	67.62 $*$
AdamW [52]	94.95	90.75	95.95	77.30	63.40	79.63	67.93 \dagger
AdaBelief [50]	95.20 \ddagger	92.25	95.71	77.02 \ddagger	68.63	77.93	70.08 \ddagger
Stagewise SGD [13]	95.23 \pm 0.19 \diamond	92.13 \pm 0.02	95.51 \pm 0.02	78.24 \pm 0.18 \diamond	69.97 \pm 0.02	78.95 \pm 0.03	70.23 \dagger
SLA [1]	95.27 \pm 0.06 \diamond	92.38 \pm 0.02	95.73 \pm 0.02	78.34 \pm 0.05 \diamond	70.20 \pm 0.04	79.54 \pm 0.02	70.30 \pm 0.09
SLRLA	95.47\pm0.20	92.63\pm0.03	96.08\pm0.07	78.58\pm0.15	70.63\pm0.02	79.85\pm0.05	70.47\pm0.12

SLRLA has better test performance than SGD and (stagewise) LA

Conclusion

- **Problems:**

(1) Why Lookahead enjoys better test performance than SGD?

LookAhead enjoys smaller excess risk error than SGD

(2) How to further improve LookAhead?

we propose a stagewise locally-regularized Lookahead with provable performance

Thanks !