

# Comp6211e: Optimization for Machine Learning

Tong Zhang

## Lecture 23: Adaptive Gradient Methods

# Stochastic Optimization in Machine Learning

In machine learning, we observe training data  $(x_i, y_i)$  for  $i = 1, \dots, n$ , and would like to learn a model parameter  $w$  of the form

$$\min_{w \in C} \left[ \frac{1}{n} \sum_{i=1}^n f_i(w) + g(w) \right].$$

More generally, we can write this optimization problem as:

$$\min_{w \in C} \phi(w), \quad \phi(w) = f(w) + g(w), \quad f(w) = \mathbf{E}_{\xi \sim D} f(\xi, w), \quad (1)$$

where  $\xi$  is a random variable, drawn from a distribution  $D$ .

# Coordinate-wise Learning Rate

Given a minibatch  $B$ , and positive definite matrix  $\Lambda$ , we define

$$f_B(w) = \frac{1}{|B|} \sum_{\xi \in B} f(\xi, w),$$

and

$$\text{prox}_{\Lambda, g}(w) = \arg \min_z \left[ \frac{1}{2} (z - w)^\top \Lambda^{-1} (z - w) + g(z) \right].$$

Consider the following general proximal stochastic gradient method below with  $B_t \sim D$ :

$$w^{(t)} = \text{prox}_{\Lambda_t, g}(w^{(t-1)} - \Lambda_t \nabla f_{B_t}(w^{(t-1)})),$$

where we replace the coordinate independent learning rate  $\eta_t$  by a diagonal matrix  $\Lambda_t = \text{diag}(\eta_{t,1}, \dots, \eta_{t,d})$ . Here  $\eta_{t,j}$  is the learning rate for the  $j$ -th coordinate at time  $t$ .

# General Property

If we define

$$Q_{\Lambda,B}(w; w') = f(w') + \nabla f_B(w')^\top (w - w') + \frac{1}{2} \|w - w'\|_{\Lambda^{-1}}^2 + g(w),$$

where

$$\|u\|_{\Lambda}^2 = u^\top \Lambda u.$$

## Proposition

*Assume that  $f(w)$  and  $g(w)$  are convex. Then*

$$f(w') + g(w) \leq Q_{\Lambda,B}(w; w') + \frac{1}{2} \|\nabla f_B(w')\|_{\Lambda}^2.$$

The first inequality follows from the same argument of Lecture 19. For the second inequality, we note that

$$\begin{aligned} f(w') + g(w) &\leq f(w') + g(w) + \frac{1}{2} \|w - w' + \Lambda \nabla f_B(w')\|_{\Lambda^{-1}}^2 \\ &= Q_{\Lambda, B}(w; w') + \frac{1}{2} \|\nabla f_B(w')\|_{\Lambda}^2. \end{aligned}$$

This proves the result.

# Convergence Result

## Theorem

Let  $\Delta_j = \max_{w, w' \in C} |w - w'|_j$ , and  $\Delta = \text{diag}(\Delta_1, \dots, \Delta_d)$ . If we take  $\Lambda_1 \geq \Lambda_2 \geq \dots$ , then

$$\mathbf{E} \sum_{t=1}^T \phi(w^{(t-1)}) \leq T\phi(w) + \mathbf{E}[g(w^{(0)}) - g(w^{(T)})] + \frac{1}{2} \mathbf{E} \text{trace}(\Delta^2 \Lambda_T^{-1}) + \frac{1}{2} \sum_{t=1}^T$$

We have from Proposition 1

$$\begin{aligned} f(\mathbf{w}^{(t-1)}) + g(\mathbf{w}^{(t)}) &\leq Q_{\Lambda_t, B_t}(\mathbf{w}^{(t)}; \mathbf{w}^{(t-1)}) + \frac{1}{2} \|\nabla f_{B_t}(\mathbf{w}^{(t-1)})\|_{\Lambda_t}^2 \\ &\leq Q_{\Lambda_t, B_t}(\mathbf{w}; \mathbf{w}^{(t-1)}) - \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^\top \Lambda_t^{-1} (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} \|\nabla f_{B_t}(\mathbf{w}^{(t-1)})\|_{\Lambda_t}^2 \\ &\leq f(\mathbf{w}^{(t-1)}) + \nabla f_{B_t}(\mathbf{w}^{(t-1)})^\top (\mathbf{w} - \mathbf{w}^{(t-1)}) + g(\mathbf{w}) \\ &\quad + \frac{1}{2} \|\mathbf{w} - \mathbf{w}^{(t-1)}\|_{\Lambda_t^{-1}}^2 - \frac{1}{2} \|\mathbf{w} - \mathbf{w}^{(t)}\|_{\Lambda_t^{-1}}^2 + \frac{1}{2} \|\nabla f_{B_t}(\mathbf{w}^{(t-1)})\|_{\Lambda_t}^2. \end{aligned}$$

We obtain

$$\mathbf{E}_{B_t}[f(\mathbf{w}^{(t-1)}) + g(\mathbf{w}^{(t)})] \leq \phi(\mathbf{w}) + \frac{1}{2} \mathbf{E}_{B_t} \left[ \|\mathbf{w} - \mathbf{w}^{(t-1)}\|_{\Lambda_t^{-1}}^2 - \|\mathbf{w} - \mathbf{w}^{(t)}\|_{\Lambda_t^{-1}}^2 \right].$$

# Proof: continues

By summing over  $t$ , and noticing that (we take  $\Lambda_0^{-1} = 0$ ):

$$\begin{aligned} & \sum_{t=1}^T \left[ \|w - w^{(t-1)}\|_{\Lambda_t^{-1}}^2 - \|w - w^{(t)}\|_{\Lambda_t^{-1}}^2 \right] \\ & \leq \sum_{t=1}^T \|w - w^{(t-1)}\|_{\Lambda_t^{-1} - \Lambda_{t-1}^{-1}}^2 \\ & \leq \sum_{t=1}^T \text{trace}(\Delta^2 (\Lambda_t^{-1} - \Lambda_{t-1}^{-1})) = \text{trace}(\Delta^2 \Lambda_T^{-1}), \end{aligned}$$

we obtain the bound.



# AdaGrad Algorithm

---

## Algorithm 1: AdaGrad

---

**Input:**  $\phi(\cdot)$ , learning rates  $\eta$ ,  $\epsilon > 0$ ,  $w^{(0)}$

**Output:**  $w^{(T)}$

```
1  $\tilde{g}_0^2 = [\epsilon, \dots, \epsilon]$ 
2 for  $t = 1, 2, \dots, T$  do
3   Randomly pick  $B_t \sim D$ 
4   Let  $g_t = \nabla_w f_{B_t}(w^{(t-1)})$ 
5   Let  $\tilde{g}_t^2 = \tilde{g}_{t-1}^2 + g_t^2$ 
6   Let  $\Lambda_t = \eta \text{diag}(\tilde{g}_t^{-1})$ 
7   Let  $w^{(t)} = \text{prox}_{\Lambda_t g}(w^{(t-1)} - \Lambda_t g_t)$ 
```

**Return:**  $w^{(T)}$

---

## Corollary

If for some  $\eta > 0$ , we take  $\eta_{t,j}$  in  $\Lambda_t$  for  $j = 1, \dots, d$  as

$$\eta_{t,j} = \frac{\eta \Delta_j}{\sqrt{\epsilon + \sum_{s=1}^t [\nabla f_{B_s}(\mathbf{w}^{(s-1)})]_j^2}},$$

then

$$\begin{aligned} \mathbf{E} \sum_{t=1}^T \phi(\mathbf{w}^{(t-1)}) &\leq T\phi(\mathbf{w}) + \mathbf{E}[g(\mathbf{w}^{(0)}) - g(\mathbf{w}^{(T)})] \\ &\quad + (0.5\eta^{-1} + \eta) \mathbf{E} \sum_{j=1}^d \Delta_j \sqrt{\epsilon + \sum_{s=1}^T [\nabla f_{B_s}(\mathbf{w}^{(s-1)})]_j^2}. \end{aligned}$$

# Proof

We have

$$\text{trace}(\Delta^2 \Lambda_T^{-1}) = \eta^{-1} \sum_{j=1}^d \Delta_j \sqrt{\epsilon + \sum_{s=1}^T [\nabla f_{B_s}(\mathbf{w}^{(s-1)})]_j^2},$$

and

$$\begin{aligned} \sum_{t=1}^T \|\nabla f_{B_t}(\mathbf{w}^{(t-1)})\|_{\Lambda_t}^2 &= \eta^2 \sum_{t=1}^T \sum_{j=1}^d \Delta_j^2 \frac{\eta_{t,j}^{-2} - \eta_{t-1,j}^{-2}}{\eta_{t,j}^{-1}} \\ &\leq 2\eta^2 \sum_{t=1}^T \sum_{j=1}^d \Delta_j^2 \frac{\eta_{t,j}^{-2} - \eta_{t-1,j}^{-2}}{\eta_{t,j}^{-1} + \eta_{t-1,j}^{-1}} = 2\eta^2 \sum_{t=1}^T \sum_{j=1}^d \Delta_j^2 (\eta_{t,j}^{-1} - \eta_{t-1,j}^{-1}) \\ &= 2\eta \text{trace}(\Delta^2 \Lambda_T^{-1}). \end{aligned}$$

This implies the bound.

---

## Algorithm 2: AdaGrad-RDA

---

**Input:**  $f(\cdot), g(\cdot), w^{(0)}, \eta_0, \eta_1, \eta_2, \dots$   
 $h(w)$  (default is  $h(w) = 0.5\eta_0\|w\|_2^2$ )

**Output:**  $w^{(T)}$

```
1 Let  $\tilde{\alpha}_0 \in \partial h(w^{(0)})$ 
2 Let  $\tilde{\Lambda}_0 = 0$ 
3  $\tilde{g}_0^2 = [\epsilon, \dots, \epsilon]$ 
4 for  $t = 1, 2, \dots, T$  do
5     Randomly select a minibatch  $B_t$  of  $m$  independent samples from  $D$ 
6     Let  $g_t = \nabla_w f_{B_t}(w^{(t-1)})$ 
7     Let  $\tilde{g}_t^2 = \tilde{g}_{t-1}^2 + g_t^2$ 
8     Let  $\Lambda_t = \eta \text{diag}(\tilde{g}_t^{-1})$ 
9     Let  $\tilde{\alpha}_t = \tilde{\alpha}_{t-1} - \Lambda_t g_t$ 
10    Let  $\tilde{\Lambda}_t = \tilde{\Lambda}_{t-1} + \Lambda_t$ 
11    Let  $w^{(t)} = \text{prox}_{\tilde{\Lambda}_t g}(\tilde{\alpha}_t)$ 
```

**Return:**  $w^{(T)}$

---

---

## Algorithm 3: RMSprop

---

**Input:**  $\phi(\cdot)$ , learning rates  $\eta$ ,  $\rho$  (default is 0.9),  $\epsilon > 0$ ,  $w^{(0)}$

**Output:**  $w^{(T)}$

```
1  $\tilde{g}_0^2 = [\epsilon, \dots, \epsilon]$ 
2 for  $t = 1, 2, \dots, T$  do
3   Randomly pick  $B_t \sim D$ 
4   Let  $g_t = \nabla_w f_{B_t}(w^{(t-1)})$ 
5   Let  $\tilde{g}_t^2 = \rho \tilde{g}_{t-1}^2 + (1 - \rho) g_t^2$ 
6   Let  $\Lambda_t = \eta \text{diag}(\tilde{g}_t^{-1})$ 
7   Let  $w^{(t)} = \text{prox}_{\Lambda_t g}(w^{(t-1)} - \Lambda_t g_t)$ 
```

**Return:**  $w^{(T)}$

---

---

## Algorithm 4: AdaDelta

---

**Input:**  $\phi(\cdot)$ , learning rates  $\eta, \rho, \epsilon > 0$ ,  $w^{(0)}$

**Output:**  $w^{(T)}$

- 1 Let  $\tilde{g}_0^2 = 0$
- 2 Let  $\eta_0^2 = 0$
- 3 **for**  $t = 1, 2, \dots, T$  **do**
- 4     Randomly pick  $B_t \sim D$
- 5     Let  $g_t = \nabla_w f_{B_t}(w^{(t-1)})$
- 6     Let  $\tilde{g}_t^2 = \rho \tilde{g}_{t-1}^2 + (1 - \rho) g_t^2$
- 7     Let  $\Lambda_t = \text{diag} \left( \sqrt{\epsilon + \eta_{t-1}^2} / \sqrt{\epsilon + \tilde{g}_t^2} \right)$
- 8     Let  $w^{(t)} = \text{prox}_{\Lambda_t g}(w^{(t-1)} - \Lambda_t g_t)$
- 9     Let  $\eta_t^2 = \rho \eta_{t-1}^2 + (1 - \rho)(w^{(t)} - w^{(t-1)})^2$

**Return:**  $w^{(T)}$

---

---

## Algorithm 5: AdaMD

---

**Input:**  $f(\cdot), g(\cdot), w^{(0)}, \eta_0, c, p$  (default is  $\lceil n/m \rceil$ )

**Output:**  $w^{(T)}$

```

1  Let  $\Lambda_0 = \eta_0 \mathbf{I}$ 
2  Let  $\tilde{g}^2 = 0$ 
3  Let  $w_{\min} = w^{(0)}$ 
4  Let  $w_{\max} = w^{(0)}$ 
5  Let  $q = 0$ 
6  for  $t = 1, 2, \dots, T$  do
7      Randomly pick  $B_t \sim D$ 
8      Let  $g_t = \nabla_{w} f_{B_t}(w^{(t-1)})$ 
9      Let  $\tilde{g}^2 = \tilde{g}^2 + g_t^2$ 
10     Let  $\Lambda_t = \Lambda_{t-1}$ 
11     Let  $w^{(t)} = \text{prox}_{\Lambda_t g}(w^{(t-1)} - \Lambda_t g_t)$ 
12     Let  $w_{\min} = \min(w_{\min}, w^{(t)})$ 
13     Let  $w_{\max} = \max(w_{\max}, w^{(t)})$ 
14     Let  $q = q + 1$ 
15     if  $q \geq p$  then
16         Let  $\Lambda_t = \text{diag}((c(w_{\max} - w_{\min}) + \sqrt{\epsilon}) / \sqrt{\epsilon + \tilde{g}^2})$ 
17         Let  $\tilde{g}^2 = 0$ 
18         Let  $w_{\min} = w^{(t)}$ 
19         Let  $w_{\max} = w^{(t)}$ 
20         Let  $q = 0$ 

```

**Return:**  $w^{(T)}$

---

We study the smoothed hinge loss function  $\phi_\gamma(z)$  with  $\gamma = 1$ , and solves the following  $L_1 - L_2$  regularization problem:

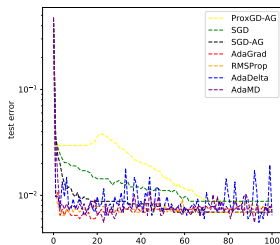
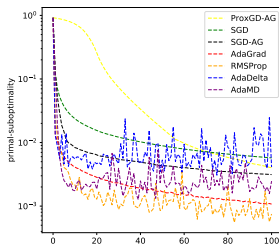
$$\min_w \left[ \underbrace{\frac{1}{n} \sum_{i=1}^n \phi_\gamma(w^\top x_i y_i)}_{f(w)} + \underbrace{\frac{\lambda}{2} \|w\|_2^2 + \mu \|w\|_1}_{g(w)} \right].$$

We compare different algorithms with constant learning rate.

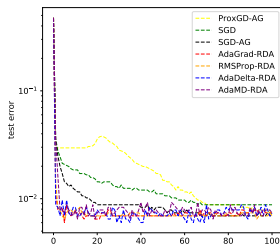
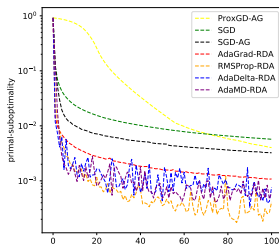


# Comparisons (proximal)

epochs ( $\gamma = 0.1$   $\lambda = 1e-05$   $\mu = 0.001$ )

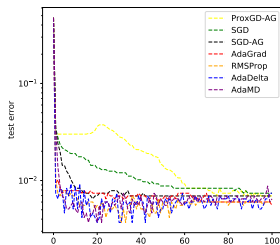
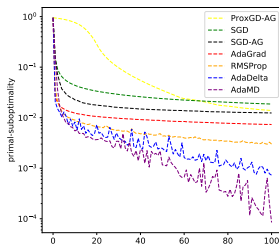


epochs ( $\gamma = 0.1$   $\lambda = 1e-05$   $\mu = 0.001$ )

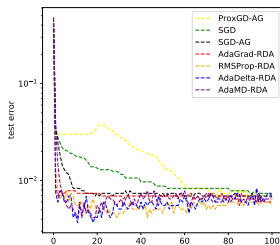
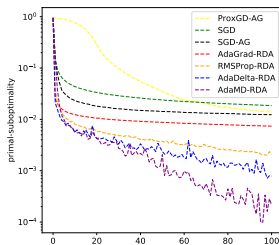


# Comparisons (no-proximal)

epochs ( $\gamma = 0.1$   $\lambda = 1e-16$   $\mu = 1e-16$ )



epochs ( $\gamma = 0.1$   $\lambda = 1e-16$   $\mu = 1e-16$ )



## Coordinate-wise Gradients

- Using gradient to estimate variance
- Only applies to stochastic optimization

## AdaGrad

- Search a larger parameter space quickly.
- $L_\infty$ -norm rather than  $L_2$ -norm
- Might hurt generalization in some cases

Faster Convergence, but less robust