# Comp6211e: Optimization for Machine Learning

Tong Zhang

Lecture 24: Applications of Adaptive Gradient Methods

# Stochastic Optimization in Machine Learning

We can write this optimization problem as:

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

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

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

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

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

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

## CTR Problem

Large scale logistic regression is useful for CTR (click through rate) estimation problem in computational advertising.

One of the standard approaches is to update models online in real time, to solve linear logistic regression:

$$\sum_{t=1}^{T} \log(1 + \exp(-w^\top x_t y_t)) + \frac{\lambda}{2}\|w\|_2^2 + \mu\|w\|_1,$$

where $y_t \in \{\pm 1\}$ indicates whether an ad is clicked or not. These are generally huge models with billions of parameters, and each feature $x_t$ is sparse.

Linear methods are still widely used, although deep learning models have been applied in various situation.

# FTRL (followed the regularized leader)

**Algorithm 1:** FTRL (follow the regularized leader)

**Input**: $f(\cdot)$, $g(\cdot)$, $w^{(0)}$, $\eta_0, \eta_1, \eta_2, \ldots$
    $h(w)$ (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 $\tilde{g}_0^2 = [\epsilon, \ldots, \epsilon]$
3 $\bar{g}_0^2 = 0$
4 **for** $t = 1, 2, \ldots, T$ **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 $\Lambda_t = \eta \mathrm{diag}(\tilde{g}_t^{-1})$
9     Let $\bar{g}_t = \bar{g}_{t-1} + g_t - (\Lambda_t^{-1} - \Lambda_{t-1}^{-1})w^{(t-1)}$
0     Let $w^{(t)} = \mathrm{prox}_{\Lambda_t g}(-\Lambda_t \bar{g}_t)$

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

# AdaGrad-RDA (from last lecture)

**Algorithm 2:** AdaGrad-RDA

**Input**: $f(\cdot)$, $g(\cdot)$, $w^{(0)}$, $\eta_0, \eta_1, \eta_2, \ldots$
$\quad\quad h(w)$ (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   $\tilde{g}_0^2 = [\epsilon, \ldots, \epsilon]$

4   **for** $t = 1, 2, \ldots, T$ **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**: $w^{(T)}$

# Original AdaGrad-RDA (from AdaGrad paper)

**Algorithm 3:** AdaGrad-RDA

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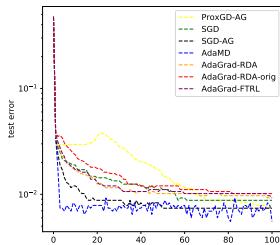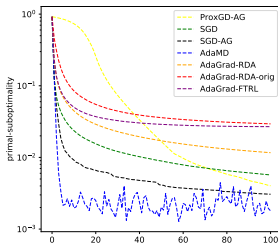   $h(w)$ (default is $h(w) = 0.5\eta_0\|w\|_2^2$)

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

**1** Let $\tilde{\alpha}_0 \in \partial h(w^{(0)})$

**2** $\tilde{g}_0^2 = [\epsilon, \ldots, \epsilon]$

**3** $\bar{g}_0^2 = 0$

**4 for** $t = 1, 2, \ldots, T$ **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 $\Lambda_t = \eta \operatorname{diag}(\tilde{g}_t^{-1})$

**9**   Let $\bar{g}_t = \bar{g}_{t-1} + g_t$

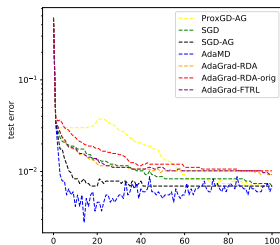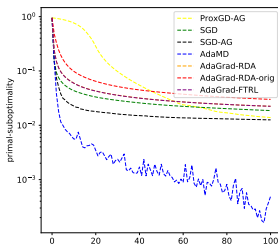**0**   Let $w^{(t)} = \operatorname{prox}_{t\Lambda_t g}(-\Lambda_t \bar{g}_t)$

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

## FTRL

**Algorithm 4:** FTRL (follow the regularized leader)

**Input**: $f(\cdot)$, $g(\cdot)$, $w^{(0)}$, $\eta_0, \eta_1, \eta_2, \ldots$
$\quad\quad$ $h(w)$ (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 $\tilde{g}_0^2 = [\epsilon, \ldots, \epsilon]$
3 $\bar{g}_0^2 = 0$
4 **for** $t = 1, 2, \ldots, T$ **do**
5 $\quad$ Randomly select a minibatch $B_t$ of $m$ independent samples from $D$
6 $\quad$ Let $g_t = \nabla_w f_{B_t}(w^{(t-1)})$
7 $\quad$ Let $\tilde{g}_t^2 = \tilde{g}_{t-1}^2 + g_t^2$
8 $\quad$ Let $\Lambda_t = \eta\,\text{diag}(\tilde{g}_t^{-1})$
9 $\quad$ Let $\bar{g}_t = \bar{g}_{t-1} + g_t - (\Lambda_t^{-1} - \Lambda_{t-1}^{-1})w^{(t-1)}$
0 $\quad$ Let $w^{(t)} = \text{prox}_{t\Lambda_t g}(-\Lambda_t\bar{g}_t)$

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

# Experiments



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

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

# AdaGrad with Acceleration

**Algorithm 4:** ADAM (without proximal term)

**Input**: $\phi(\cdot)$, learning rates $\eta$ , $\beta$ (default is 0.9), $\rho$ (default is 0.999),
$\quad\quad\epsilon > 0$ (default is $10^{-8}$ , $w^{(0)}$

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

1 Let $\tilde{g}_0^2 = 0$
2 Let $\bar{g}_0 = 0$
3 **for** $t = 1, 2, \ldots, T$ **do**
4 $\quad$ Randomly pick $B_t \sim D$
5 $\quad$ Let $g_t = \nabla_w f_{B_t}(w^{(t-1)})$
6 $\quad$ Let $\tilde{g}_t^2 = \rho \tilde{g}_{t-1}^2 + (1 - \rho)g_t^2$
7 $\quad$ Let $\bar{g}_t = \beta \bar{g}_{t-1} + g_t$
8 $\quad$ Let $\Lambda_t = \eta \text{diag}((\tilde{g}_t/(1 - \rho) + \epsilon)^{-1})$
9 $\quad$ Let $w^{(t)} = w^{(t-1)} - \Lambda_t \bar{g}_t$

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

## Proximal Version

**Algorithm 5:** Prox-ADAM

**Input**: $\phi(\cdot)$, learning rates $\eta$ , $\beta$ (default is 0.9), $\rho$ (default is 0.999),
$\epsilon > 0$ (default is $10^{-8}$ , $w^{(0)}$

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

1 Let $\tilde{g}_0^2 = 0$
2 Let $\bar{g}_0 = 0$
3 **for** $t = 1, 2, \ldots, T$ **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 $\bar{g}_t = \beta \bar{g}_{t-1} + g_t$
8     Let $\Lambda_t = \eta \mathrm{diag}((\tilde{g}_t + \epsilon)^{-1})$
9     Let $w^{(t)} = \mathrm{prox}_{\Lambda_t g}(w^{(t-1)} - \Lambda_t \bar{g}_t)$
0     $\bar{g}_t = \Lambda_t^{-1}(w^{(t-1)} - w^{(t)})$

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

## Heavy-Ball Version

**Algorithm 6:** Prox-AdaHB

**Input**: $\phi(\cdot)$, learning rates $\eta$, $\beta$ (default is 0.9), $\rho$ (default is 0.999), $\epsilon > 0$ (default is $10^{-8}$, $w^{(0)}$

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

1 Let $\tilde{g}_0^2 = 0$

2 Let $p_0 = 0$

3 **for** $t = 1, 2, \ldots, T$ **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 $\Lambda_t = \eta \mathrm{diag}((\tilde{g}_t + \epsilon)^{-1})$

8      Let $w^{(t)} = \mathrm{prox}_{\Lambda_t g}(w^{(t-1)} + \beta p_{t-1} - \Lambda_t g_t)$

9      Let $p_t = w^{(t)} - w^{(t-1)}$

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

## Nesterov's Acceleration

**Algorithm 7:** Prox-AdaACCL

**Input**: $\phi(\cdot)$, learning rates $\eta$, $\beta$ (default is 0.9), $\rho$ (default is 0.999), $\epsilon > 0$ (default is $10^{-8}$, $w^{(0)}$

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

1   $\tilde{g}_0^2 = 0$

2   $p_0 = 0$

3   **for** $t = 1, 2, \ldots, T$ **do**

4      Randomly pick $B_t \sim D$

5      Let $v^{(t)} = w^{(t-1)} + \beta(w^{(t-1)} - w^{(t-2)})$

6      Let $g_t = \nabla_w f_{B_t}(v^{(t-1)})$

7      Let $\tilde{g}_t^2 = \rho \tilde{g}_{t-1}^2 + (1 - \rho)g_t^2$

8      Let $\Lambda_t = \eta \mathrm{diag}((\tilde{g}_t + \epsilon)^{-1})$

9      Let $w^{(t)} = \mathrm{prox}_{\Lambda_t g}(v^{(t)} - \Lambda_t g_t)$

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

# RDA version

## Algorithm 8: Prox-AdaACCL-RDA

**Input**: $\phi(\cdot)$, learning rates $\eta$, $\beta$ (default is 0.9), $\rho$ (default is 0.999), $\epsilon > 0$ (default is $10^{-8}$, $w^{(0)}$
**Output**: $w^{(T)}$

1   Let $\tilde{g}_0^2 = 0$
2   Let $u^{(0)} = w^{(0)}$
3   Let $\tilde{\alpha}_0 = 0$
4   Let $\tilde{\Lambda}_0 = 0$
5   **for** $t = 1, 2, \ldots, T$ **do**
6      Randomly pick $B_t \sim D$
7      Let $v^{(t)} = \beta w^{(t-1)} + (1 - \beta)u^{(t-1)}$
8      Let $g_t = \nabla_w f_{B_t}(v^{(t)})$
9      Let $\tilde{g}_t^2 = \rho\tilde{g}_{t-1}^2 + (1 - \rho)g_t^2$
10      Let $\Lambda_t = \frac{\eta}{1-\beta}\mathrm{diag}((\tilde{g}_t + \epsilon)^{-1})$
11      Let $\tilde{\alpha}_t = \tilde{\alpha}_{t-1} - \Lambda_t g_t$
12      Let $\tilde{\Lambda}_t = \tilde{\Lambda}_{t-1} + \Lambda_t$
13      Let $u^{(t)} = \mathrm{prox}_{\tilde{\Lambda}_t g}(\tilde{\alpha}_t)$
14      Let $w^{(t)} = \beta w^{(t-1)} + (1 - \beta)u^{(t)}$

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

## **Algorithm 9:** Prox-AdaACCL-v3

**Input**: $\phi(\cdot)$, learning rates $\eta$ , $\beta$ (default is 0.9), $\rho$ (default is 0.999), $\epsilon > 0$ (default is $10^{-8}$ , $w^{(0)}$

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

1 Let $\tilde{g}_0^2 = 0$

2 Let $u^{(0)} = w^{(0)}$

3 Let $\tilde{\alpha}_0 = 0$

4 Let $\tilde{\Lambda}_0 = 0$

5 **for** $t = 1, 2, \ldots, T$ **do**

6      Randomly pick $B_t \sim D$

7      Let $v^{(t)} = \beta w^{(t-1)} + (1 - \beta) u^{(t-1)}$

8      Let $g_t = \nabla_w f_{B_t}(v^{(t)})$

9      Let $\tilde{g}_t^2 = \rho \tilde{g}_{t-1}^2 + (1 - \rho) g_t^2$

10      Let $\Lambda_t = \frac{\eta}{1-\beta} \operatorname{diag}((\tilde{g}_t + \epsilon)^{-1})$

11      Let $u^{(t)} = \operatorname{prox}_{\Lambda_t g}(u^{(t-1)} - \Lambda_t g_t)$

12      Let $w^{(t)} = \beta w^{(t-1)} + (1 - \beta) u^{(t)}$

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

# Empirical Studies

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

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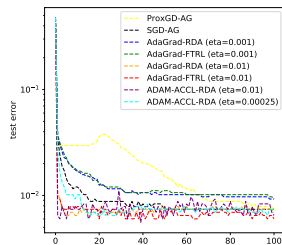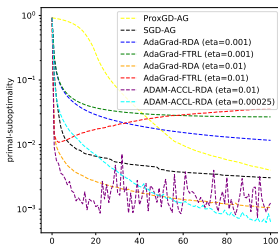
We compare different algorithms with constant learning rate.

# Comparisons (proximal)
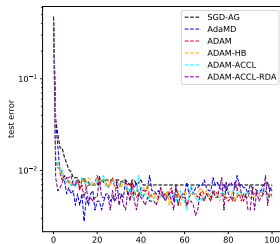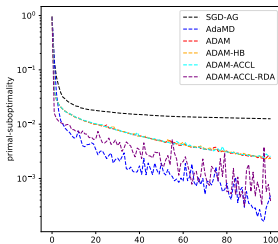


epochs ($\gamma = 0.1$ $\lambda$=1e-05 $\mu$=0.001)
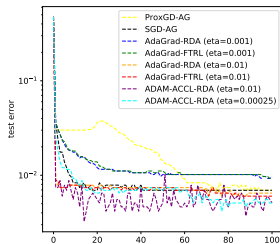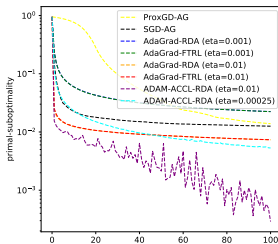


epochs ($\gamma = 0.1$ $\lambda$=1e-05 $\mu$=0.001)

# Comparisons (no-proximal)

# Summary

This course introduced basic concepts used in optimization for machine learning

- First order gradient methods
- Acceleration
- Primal dual methods
- Coordinate descent
- Proximal methods
- Dual averaging
- Variance reduction
- Coordinate-wise Learning rate (AdaGrad)
- ...

By correctly combining these concepts, we can get very effective algorithms for machine learning.