

# Comp6211e: Optimization for Machine Learning

Tong Zhang

## Lecture 21: Stochastic Adaptive Learning Rate and Acceleration 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  $D$ .

# Gradient versus Stochastic Gradient

In gradient based methods, we use gradient

$$\nabla f(\boldsymbol{w}) = \mathbf{E}_{\xi} \nabla f(\xi, \boldsymbol{w}).$$

In SGD, we replace the full gradient with stochastic gradient

$$\nabla_{\boldsymbol{w}} f(\xi, \boldsymbol{w}^{(t-1)}),$$

or minibatch stochastic gradient:

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

# Adaptive Learning Rate

Consider a minibatch  $B$ , and let

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

and let

$$V = \mathbf{E}_{\xi} \|\nabla_{\mathbf{w}} f(\xi, \mathbf{w}) - \nabla_{\mathbf{w}} f(\mathbf{w})\|_2^2.$$

Now consider the SGD update rule

$$\mathbf{w}' = \mathbf{w} - \eta \nabla f_B(\mathbf{w}),$$

How to determine the best learning rate  $\eta$ ?

# One-step Objective Reduction

$$\begin{aligned}\mathbf{E}_B f(w') &= \mathbf{E}_B f(w - \eta \nabla f_B(w)) \\ &\leq \mathbf{E}_B \left[ f(w) - \eta \nabla f(w)^\top \nabla f_B(w) + \frac{\eta^2 L}{2} \mathbf{E}_B \|\nabla f_B(w)\|_2^2 \right] \\ &= f(w) - (\eta - 0.5\eta^2 L) \|\nabla f(w)\|_2^2 + \frac{\eta^2 L}{2m} V.\end{aligned}$$

If we choose

$$\eta \leq \frac{\|\nabla f(w)\|_2^2}{L(\|\nabla f(w)\|_2^2 + V/m)},$$

then we have convergence rate of

$$\mathbf{E}_B f(w') = f(w) - 0.5\eta \|\nabla f(w)\|_2^2. \quad (2)$$

If

$$V \leq (n - m) \|\nabla f(\mathbf{w})\|_2^2, \quad (3)$$

then (2) holds as long as

$$\eta \leq m/(nL).$$

Since each minibatch SGD requires  $m$  gradient computations, per sample error reduction as in (2) becomes

$$\frac{\eta}{m} \|\nabla f(\mathbf{w})\|_2^2 \geq \frac{\|\nabla f(\mathbf{w})\|_2^2}{nL},$$

which is the per sample function value reduction as GD with step size  $\eta = 1/L$ :

$$\frac{\eta}{n} \|\nabla f(\mathbf{w})\|_2^2 = \frac{\|\nabla f(\mathbf{w})\|_2^2}{nL}.$$

# Generalized Armijo-Goldstein Condition

## Proposition (Stochastic AG Criterion)

Consider the SGD update rule:

$$w' = w - \eta \nabla f_B(w).$$

Let

$$V_{B',B} = \left[ f_{B'}(w - \eta \nabla f_B(w)) - \eta \nabla f_{B'}(w)^\top \nabla f_B(w) - f_{B'}(w) \right].$$

If for some  $c < 1$ ,

$$\mathbf{E}_{B',B} V_{B',B} \leq c\eta \|\nabla f(w)\|_2^2,$$

then

$$\mathbf{E}_B f(w') \leq f(w) - (1 - c)\eta \|\nabla f(w)\|_2^2.$$

When  $g(w) \neq 0$ , we may generalize the Proposition 1 by using gradient mapping as follows.

$$\begin{aligned}\text{prox}_{\eta g}(w) &= \arg \min_z \left[ \frac{1}{2\eta} \|z - w\|_2^2 + g(z) \right] \\ D_{\eta} \phi(w) &= \frac{1}{\eta} (w - \text{prox}_{\eta g}(w - \eta \nabla f(w))) \\ D_{\eta, B} \phi(w) &= \frac{1}{\eta} (w - \text{prox}_{\eta g}(w - \eta \nabla f_B(w))).\end{aligned}$$

With this modification, we obtain the adaptive learning rate method for stochastic gradient descent.



# Stochastic Gradient with AG Learning Rate

---

## Algorithm 1: Stochastic Proximal Gradient Descent with Adaptive Learning Rate

---

**Input:**  $f(\cdot), g(\cdot), w_0, \eta_0, \rho$  (default is  $\lceil n/m \rceil$ )

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

```
1 Let  $\eta_1 = \eta_0$ 
2 Let  $q = 0$ 
3 Let  $V = 0$ 
4 Let Randomly select a minibatch  $B_0$  independent samples from  $D$ 
5 Let  $\tilde{g} = D_{\eta_0, B_0} \phi(w_0)$ 
6 for  $t = 1, 2, \dots, T$  do
7     Randomly select a minibatch  $B_t$  of  $m$  independent samples from  $D$ 
8     Let  $w^{(t)} = \text{prox}_{\eta_t g}(w^{(t-1)} - \eta_t \nabla f_{B_t}(w^{(t-1)}))$ 
9     Let  $V = V + [f_{B_t}(w^{(t)} - \eta_t \tilde{g}) - f_{B_t}(w^{(t)}) - \eta_t \nabla f_{B_t}(w^{(t)})^\top \tilde{g}]$ 
10    Let  $\tilde{g} = (w^{(t-1)} - w^{(t)}) / \eta_t$ 
11    Let  $\eta_{t+1} = \eta_t$ 
12    Let  $q = q + 1$ 
13    if  $q \geq \rho$  then
14        Let  $\rho = \min(0.5, \max(2, qc\eta_t \|D_{\eta_t} \phi(w_t)\|_2^2 / V))$ 
15        Let  $\eta_{t+1} = \eta_t \rho$ 
16        Let  $q = 0$ 
17        Let  $V = 0$ 
18        Let  $\tilde{g} = D_{\eta_{t+1}, B_t} \phi(w_t)$ 
```

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

---

# Stochastic Accelerated Gradient

---

**Algorithm 2:** Stochastic Accelerated Proximal Gradient Descent

---

**Input:**  $f(\cdot)$ ,  $g(\cdot)$ ,  $\{\eta_t, \beta_t\}$

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

```
1 for  $t = 1, 2, \dots, T$  do
2   Let  $z^{(t)} = w^{(t-1)} + \beta_t(w^{(t-1)} - w^{(t-2)})$ 
3   Randomly select a minibatch  $B_t$  of  $m$  independent samples from  $D$ 
4   Let  $w^{(t)} = \text{prox}_{\eta_t g}(z^{(t)} - \eta_t \nabla f_{B_t}(z^{(t)}))$ 
```

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

---

# Acceleration with AG Learning Rate

## Algorithm 3: Stochastic Accelerated Proximal Gradient with Adaptive Learning Rate

**Input:**  $f(\cdot), g(\cdot), w_0, \eta_0, \beta, p$  (default is  $\lceil n/m \rceil$ )

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

```
1 Let  $\eta_1 = \eta_0$ 
2 Let  $q = 0$ 
3 Let  $V = 0$ 
4 Let Randomly select a minibatch  $B_0$  independent samples from  $D$ 
5 Let  $\tilde{g} = D_{\eta_0, B_0} \phi(w_0)$ 
6 for  $t = 1, 2, \dots, T$  do
7     Let  $z^{(t)} = w^{(t-1)} + \beta(w^{(t-1)} - w^{(t-2)})$ 
8     Randomly select a minibatch  $B_t$  of  $m$  independent samples from  $D$ 
9     Let  $w^{(t)} = \text{prox}_{\eta_t g}(z^{(t)} - \eta_t \nabla f_{B_t}(z^{(t)}))$ 
10    Let  $V = V + [f_{B_t}(z^{(t)} - \eta_t \tilde{g}) - f_{B_t}(z^{(t)}) - \eta_t \nabla f_{B_t}(z^{(t)})^\top \tilde{g}]$ 
11    Let  $\tilde{g} = (z^{(t)} - w^{(t)}) / \eta_t$ 
12    Let  $\eta_{t+1} = \eta_t$ 
13    Let  $q = q + 1$ 
14    if  $q \geq p$  then
15        Let  $\rho = \max(0.5, \min(2, qc(1 - \beta)\eta_t \|D_{\eta_t} \phi(w_t)\|_2^2 / V))$ 
16        Let  $\eta_{t+1} = \eta_t \rho$ 
17        Let  $q = 0$ 
18        Let  $V = 0$ 
19        Let  $\tilde{g} = D_{\eta_{t+1}, B_t} \phi(w_t)$ 
```

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

# Momentum (Heavy Ball) Method

---

**Algorithm 4:** Stochastic Heavy-Ball Gradient Descent

---

**Input:**  $f(\cdot)$ ,  $g(\cdot)$ ,  $\{\eta_t, \beta_t\}$

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

```
1 for  $t = 1, 2, \dots, T$  do
2   Let  $z^{(t)} = w^{(t-1)} + \beta_t(w^{(t-1)} - w^{(t-2)})$ 
3   Randomly select a minibatch  $B_t$  of  $m$  independent samples from  $D$ 
4   Let  $w^{(t)} = \text{prox}_{\eta_t g}(z^{(t)} - \eta_t \nabla f_{B_t}(w^{(t-1)}))$ 
```

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

---

---

**Algorithm 5:** Stochastic Accelerated Regularized Dual Averaging

---

**Input:**  $f(\cdot)$ ,  $g(\cdot)$ ,  $w^{(0)}$ ,  $\tilde{\eta}_0 \leq \tilde{\eta}_1, \dots, \theta_t$   
 $h(w)$  (default is  $h(w) = 0.5\|w\|_2^2$ )

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

```
1 Let  $\tilde{\alpha}_0 \in \partial h(w^{(0)})$ 
2 Let  $v^{(0)} = w^{(0)}$ 
3 for  $t = 1, 2, \dots, T$  do
4   Let  $u^{(t)} = (1 - \theta_t)w^{(t-1)} + \theta_tv^{(t-1)}$ 
5   Randomly select a minibatch  $B_t$  of  $m$  independent samples from  $D$ 
6   Let  $\tilde{\alpha}_t = (1 - \theta_t)\tilde{\alpha}_{t-1} - \theta_t\nabla f_B(u^{(t)})$ 
7   Let  $v^{(t)} = \arg \min_w \left[ -\tilde{\alpha}_t^\top w + \tilde{\eta}_t^{-1} h(w) + g(w) \right]$ 
8   Let  $w^{(t)} = (1 - \theta_t)w^{(t-1)} + \theta_tv^{(t)}$ 
```

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

---

# Stochastic Accelerated Linearized ADMM

---

**Algorithm 6:** Stochastic Accelerated Linearized ADMM

---

**Input:**  $\phi(\cdot)$ ,  $A$ ,  $B$ ,  $c$ ,  $\{\eta_t\}$ ,  $\beta$ ,  $\rho$ ,  $\alpha_0$ ,  $w_0$ ,  $z_0$

**Output:**  $w_T, z_T, \alpha_T$

```
1 Let  $\bar{w}_0 = x_0$ 
2 Let  $\bar{z}_0 = z_0$ 
3 for  $t = 1, 2, \dots, T$  do
4   Let  $\tilde{z}_t = \bar{z}_{t-1} - \eta_t B^\top [\alpha_{t-1} + \rho(A\bar{w}_{t-1} + Bz_{t-1} - c)]$ 
5   Let  $z_t = \arg \min_z [0.5\|z - \tilde{z}_t\|_2^2 + \eta_t g(z)]$ 
6   Let  $w_t = \bar{w}_{t-1} - \eta_t \nabla f_B(\bar{w}_{t-1}) - \eta_t A^\top [\alpha_{t-1} + \rho(A\bar{w}_{t-1} + Bz_t - c)]$ 
7   Let  $\alpha_t = \alpha_{t-1} + \rho(1 - \beta)[Aw_t + Bz_t - c]$ 
8   Let  $\bar{z}_t = z_t + \beta(z_t - z_{t-1})$ 
9   Let  $\bar{w}_t = w_t + \beta(w_t - w_{t-1})$ 
```

**Return:**  $w_T, z_T, \alpha_T$

---

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 (strongly convex)

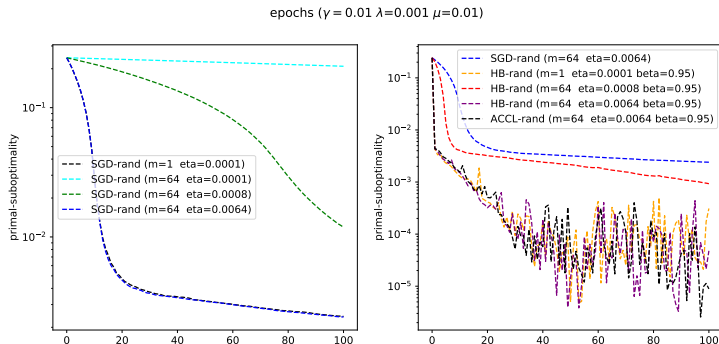


Figure: Comparisons of stochastic algorithms



# Comparisons (strongly convex)

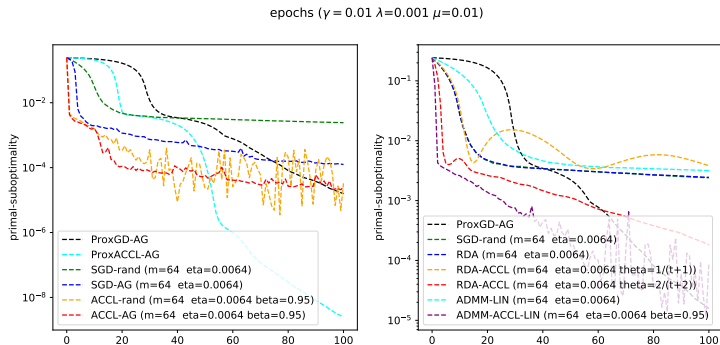
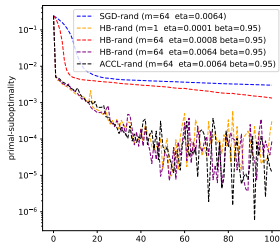
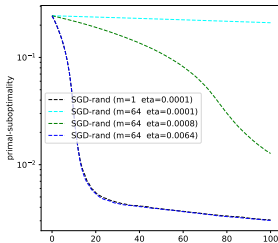


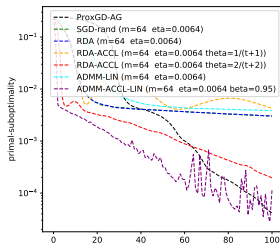
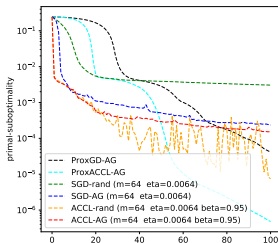
Figure: Comparisons of stochastic algorithms

# Comparisons (near non-strongly convex)

epochs ( $\gamma = 0.01$   $\lambda = 1e-05$   $\mu = 0.01$ )



epochs ( $\gamma = 0.01$   $\lambda = 1e-05$   $\mu = 0.01$ )



## SGD versus GD

- SGD converges faster when gradient is relatively large

$$\|D_{\eta}\phi(w)\|_2^2 \geq O(m/n) \times \text{minibatch variance}.$$

- Eventually GD converges faster

## Tuning of Learning Rate

- Can generalize AG criterion to the stochastic setting
- Requires estimating the variance on the fly.

## Acceleration

- improve over non-accelerated in the beginning
- eventually variance dominate the convergence rate
- has an effect of using a larger learning rate