

Learning Theory

X - domain $\in \mathbb{R}^d$

$D \sim X \times Y$

Y - label set

l - loss $f: \mathcal{H} \times X \times Y$

\mathcal{H} - hypothesis class \rightarrow parameterized by w
 $\forall w$ s.t. $\|w\|^2 \leq B$

Goal

Output: $\operatorname{argmin}_{h \in \mathcal{H}} E_{(x,y) \sim D} [l(h, x, y)]$

Algorithm

initialize w_1

For $t = 1, 2, \dots$

sample $z \sim D$

$v_t \in \partial l(w_t, z)$

$w_{t+1} = w_t - \eta_t v_t$

Output $\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$

Dr:

$S: ((x_1, y_1), \dots, (x_n, y_n))$

Output $\operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n l(h, x_i, y_i)$

Algorithm: SG

initialize w_1

For $t = 1, 2, \dots$

sample $i_t \sim \mathcal{U}([n])$

$v_t \in \partial l(w_t, x_{i_t}, y_{i_t})$

$w_{t+1} = w_t - \eta_t v_t$

Output $\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$

Denote $f_i(w) = l(w, x_i, y_i)$

$g(w) = \frac{1}{n} \sum f_i(w)$

$g(w)$ is ℓ -lipschitz

$\forall i, \forall u, v \in \mathbb{R}^d$

$\|f_i(u) - f_i(v)\| \leq \ell \|u - v\|$

$g(w)$ is L -smooth

$\forall i, \forall u, v \in \mathbb{R}^d$

$\|f'_i(u) - f'_i(v)\| \leq L \|u - v\|$

or

$\forall w, |\text{eigenvalues}[g''(w)]| \leq L$

FC

$v_t = \frac{1}{n} \sum_{i=1}^n \partial l(w_t, x_i, y_i)$

$w_{t+1} = w_t - \eta_t v_t$

$g(w)$ λ -strongly convex

$$\kappa = \frac{L}{\mu}$$

$\forall w$ eigenvalues $[g''(w)] \geq \lambda$

G

ϵ -error

complexity

convex

$$g(w^t) - g(w^*) \leq O(1/t) \rightarrow O(1/\epsilon) \quad \left[O(nd) \right]$$

strong
convexity

$$g(w^t) - g(w^*) \leq O((1 - \lambda/L)^t) = O(e^{-t} \lambda/L) \rightarrow O(\kappa \log \frac{1}{\epsilon})$$

Step size backtracking line search

G

convex

$$E[g(w^t)] - g(w^*) \leq O(1/\sqrt{t}) \rightarrow O(1/\epsilon^2)$$

λ -strong
convexity

$$E[g(w^t)] - g(w^*) \leq O(1/t) \rightarrow O(1/\epsilon)$$

Step size

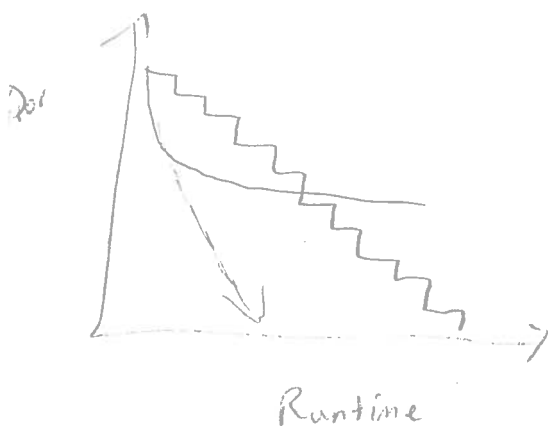
complexity

$$\eta_t = \frac{1}{L\sqrt{t}}$$

$$O(d)$$

$$\eta_t = \frac{1}{\lambda t}$$

$$O(d)$$



SAG (Stochastic Average Gradient)

• Store $\nabla^T f_i(w) \quad \forall i$

• Sample $i_t \sim \mathcal{U}([n])$

• $w_t = w_{t-1} - \eta_t \frac{1}{n} \sum \gamma_i^t$ where $\gamma_i^t = \begin{cases} \nabla f(w_{t-1}) & i = i_t \\ \gamma_i^{t-1} & \text{o.w.} \end{cases}$

initialize $\text{sum} = 0, \gamma_i = 0 \quad \forall i$

for $t = 1, 2, \dots$

$i_t \sim \mathcal{U}([n])$

$\text{sum} = \text{sum} - \gamma_{i_t} + \nabla f_{i_t}(w^t)$

$\gamma_{i_t} = \nabla f_{i_t}(w^t)$

$w^{t+1} = w^t - \frac{\eta}{n} \text{sum}$

Thm

SAG

Step size

$$\eta = \frac{1}{16L}$$

convex

$$\mathbb{E}[g(w^t)] - g(w^*) \leq O(1/t)$$

λ -strong convexity

$$\eta = \frac{1}{16L}$$

$$\mathbb{E}[g(w^t)] - g(w^*) \leq O\left(\left(1 - \min\left(\frac{\lambda}{16L}, \frac{1}{8n}\right)\right)^t\right)$$

Complexity

$$O(d)$$

$$O(d)$$

Storage

$$O(nd)$$

$$O(nd)$$

Variance Reduction

Reduce variance of sampled x by sampling y
with known expectation

$$Z_\alpha = \alpha(x - y) + E[y]$$

$$E[Z_\alpha] = \alpha E[x] + (1 - \alpha)E[y]$$

$$\text{Var}(Z_\alpha) = \alpha^2 [\text{var}(x) + \text{var}(y) - 2\text{cov}(x, y)]$$

