# Optimization Algorithms for Machine/Deep Learning

## Deterministic optimization methods

$min f(x), \quad x \in X$

$f$: differentiable

$||\nabla f(x) - \nabla f(y)|| \leq L||x-y||, \quad \forall x, y \in X$

e.g. $f(x) = \frac{1}{2}X^{\mathrm{T}}AX - b^{\mathrm{T}}x$

$\qquad \nabla f(x) = Ax - b$

$\qquad ||\nabla f(x) - \nabla f(y)|| = ||A(x-y)|| \leq ||A||||x-y||, \quad L = ||A||$

**Lemma** $\quad f(x) \leq f(y) + <\nabla f(y), x-y> + \frac{1}{2}||x-y||^2$

Proof: $\quad$ Let $\phi(t) = f(y + t(x-y))$

$\qquad\qquad \nabla \phi(t) = \nabla f(y + t(x-y))^{\mathrm{T}}(x-y)$

$\qquad\qquad \phi(1) - \phi(0) = \int_0^1 \nabla \phi(t) dt$

$\qquad\qquad f(x) - f(y) = \int_0^1 \nabla f(y + t(x-y))^{\mathrm{T}}(x-y) dt$

$f(x) - f(y) - <\nabla f(y), x-y> = \int_0^1 \nabla f(y + t(x-y))^{\mathrm{T}}(x-y) dt - \int_0^1 \nabla f(y)^{\mathrm{T}}(x-y) dt$

$\qquad\qquad = \int_0^1 (\nabla f(y + t(x-y)) - \nabla f(y))^{\mathrm{T}}(x-y) dt$

$\qquad\qquad \leq \int_0^1 ||\nabla f(y + t(x-y))|| ||x-y|| dt \quad$ (Cauchy inequality)

$\qquad\qquad \leq tL||x-y||^2 dt$

$\qquad\qquad = \frac{L}{2}||x-y||^2$

$x_{t+1} = argmin_{x \in X} Y_t <\nabla f(x), x> + \frac{1}{2}||x - x_t||^2$

**Optimality condition for the above subproblem**

- $<\gamma_t \nabla f(x_t) + x_{t+1} - x_t, x - x_{t+1}> \geq 0, \quad \forall x \in X \qquad$ (OPT1)
- $\gamma <\nabla f(x_t), x_{t+1} - x> \leq \frac{1}{2}||x - x_t||^2 - \frac{1}{2}||x - x_{t+1}|| - \frac{1}{2}||x_t - x_{t+1}||^2 \qquad$ (OPT2)

## Observation

$f(x_{t+1}) \leq f(x_t)$ if $\gamma_t \leq \frac{2}{L}$

Fix $x = x_t$ in OPT1, $<\nabla f(x_t), x_{t+1} - x_t> \leq -\frac{1}{\gamma_t}||x_{t+1} - x_t||^2$

Also, $f(x_{t+1}) \leq f(x_t) + <\nabla f(x_t), x_{t+1} - x_t> + \frac{L}{2}||x_{t+1} - x_t||^2$

$$\leq f(x_t) - (\tfrac{1}{\gamma_t} - \tfrac{L}{2})\|x_{t+1} - x_t\|^2$$

$$\leq f(x_t)$$

$$f(x_{t+1}) \leq f(x_t) + <\nabla f(x_t), x_{t+1} - x_t> + \tfrac{1}{2}\|x_{t+1} - x\|^2$$

$$= f(x_t) + <\nabla f(x_t), x - x_t> + <\nabla f(x_t), x_{t+1} - x> + \tfrac{L}{2}\|x_{t+1} - x_t\|^2$$

(Strong) Convexity          OPT2

$$\leq f(x) + \tfrac{1}{2\gamma_t}[\|x - x_t\|^2 - \|x - x_{t+1}\|^2] - \tfrac{1}{2}(\tfrac{1}{\gamma_t} - L)\|x_{t+1} - x_t\|^2$$

$$\leq f(x) + \tfrac{1}{2\gamma_t}[\|x - x_t\|^2 - \|x - x_{t+1}\|^2]$$

If $\gamma_t = \gamma \leq \tfrac{1}{L}$, $t \leq 1, 2, \ldots$, then $\sum_{t=1}^{k}[f(x_{t+1}) - f(x)] \leq \tfrac{1}{2\gamma}[\|x - x_t\|^2 - \|x - x_{t+1}\|^2]$

Notice $f(x_t) \geq f(x_{k+1})$, $\forall t \leq k+1$

$$\sum_{t=1}^{k}[f(x_{t+1}) - f(x)] \geq k[f(x_{k+1}) - f(x)]$$

Then $f(x_{k+1}) - f(x) \leq \tfrac{1}{2\gamma k}\|x - x_1\|^2$

$$\gamma = \tfrac{1}{L} \Rightarrow f(x_{k+1}) - f(x^*) \leq \tfrac{L}{2k}\|x^* - x_1\|^2$$

## Strong Convexity

$$f(x) \geq f(y) + <\nabla f(y), x - y> + \tfrac{\mu}{2}\|x - y\|^2$$

e.g. $f(x) = \tfrac{1}{2}x^{\mathrm{T}}Ax + b^{\mathrm{T}}x$, $\mu = \lambda_{min}(A)$

$$f(x_{t+1}) \leq f(x) - \tfrac{\mu}{2}\|x - x_t\|^2 + \tfrac{1}{2\gamma}[\|x - x_t\|^2 - \|x - x_{t+1}\|^2]$$

$$f(x+1) - f(x^*) + \tfrac{1}{2\gamma}\|x^* - x_{t+1}\|^2 \leq \tfrac{1}{2}(\tfrac{1}{\gamma} - \mu)\|x^* - x_t\|^2$$

$$\|x_{t+1} - x^*\| \leq (1 - \gamma\mu)\|x^* - x_t\|^2$$

If $r = \tfrac{1}{L}$, then $\|x_{t+1} - x^*\| \leq (1 - \tfrac{\mu}{L})\|x_t - x^*\|^2$

$$\|x_{k+1} - x^*\| \leq (1 - \tfrac{\mu}{L})^k\|x_t - x^*\|^2$$

If we want to have $\|x_{k+1} - x^*\| \leq \varepsilon$

It suffices to have $(1 - \tfrac{\mu}{L})^k\|x_1 - x^*\|^2 \leq \varepsilon$

$$(1 - \tfrac{\mu}{L})^k \leq \tfrac{\varepsilon}{\|x_1 - x^*\|^2}$$

$$k \cdot log(1 - \tfrac{\mu}{L}) \leq log\tfrac{\varepsilon}{\|x_i - x^*\|^2}$$

$$k \cdot (-log(1 - \tfrac{\mu}{L})) \geq log\tfrac{\|x_i - x^*\|^2}{\varepsilon}$$

$$k \geq \tfrac{1}{-log(1-\tfrac{\mu}{L})}log\tfrac{\|x_1 - x^*\|^2}{\varepsilon} \Leftarrow l \geq \tfrac{L}{\mu}log\tfrac{\|x_i - x^*\|^2}{\varepsilon}$$

$\varepsilon$: conditional number

$$\nabla f(x)$$

$E[G(x_t, \xi_t)] = \nabla f(x_t)$

Define $\delta_t = \nabla f(x_t) - G(x_t, \xi_t)$

- $E[\delta_t] = 0$

  $\delta_t$ independent of $x_t$

- $E[||\delta_t||^2] \le \sigma^2$

$x_{t+1} = argmin_{x \in X} \gamma_t < G(x_t, \xi_t), x > + \frac{1}{2}||x - x_t||^2$

$\gamma_t < G(x_t, \xi_t), x_{t+1} - x > \le \frac{1}{2}[||x - x_t||^2 - ||x - x_{t+1}||^2 - ||x_t - x_{t+1}||^2]$  (OPT2')

`todo:`

Will be inplemented later or you can pull requests my [Github Repo](#)

## Comments

1. $f(x) = E_\xi[F(x, \xi)]$, $\xi$ is continuous random variable, SGD nearly optimal
2. $f(x) = \frac{1}{N}\sum_{t=1}^{N} f_i(x)$, using randomized incremental gradient method, we can improve the speed of convergence in terms of the dependence on $\varepsilon$. But the convergence depends on $N$.

- Deep Learning

- Burer-Monteiro Law Rank

  decomposition

  $X = LU, \;\; L \in \mathbb{R}^{m \times r}, U \in \mathbb{R}^{r \times n}$

  $min_{L,U}||X - LU||^2$

## Nonconvex Optimization

$min_{x \in \mathbb{R}^n} f(x)$

$f$ is smooth but not necessarily convex

$||\nabla f(x) - \nabla f(y)|| \le L||x - y||, \;\; \forall x, y$

$x_{t+1} = x_t - \gamma_t \nabla f(x_t)$

$f(x_{t+1}) \le f(x_t) + \nabla f(x_t)^{\mathrm{T}}(x_{t+1} - x_t) + \frac{L}{2}||x_{t+1} - x_t||^2$

$\qquad = f(x_1) - \gamma_t ||\nabla f(x_t)||^2 + \frac{L\gamma_t^2}{2}||\nabla f(x_t)||^2$

$\qquad = f(x_t) - \gamma_t(1 - \frac{L\gamma_t}{2})||\nabla f(x_t)||^2$

$\gamma_t(1 - \frac{L\gamma_t}{2})||\nabla f(x_t)||^2 \le f(x_t) - f(x_{t+1})$

$\sum_{t=1}^{k} \gamma_t(1 - \frac{L\gamma_t}{2})||\nabla f(x_t)||^2 \le f(x_1) - f(x_{t+1}) \le f(x_1) - f^*$

Output $\overline{x_k}$ s.t. $||\nabla f(\overline{x_k})|| = min_{t=1,\ldots,k}||\nabla f(\overline{x_t})||$

$0 \le \gamma_t \le \frac{2}{L}$

$$\sum_{t=1}^{k} \gamma_t (1 - \frac{L\gamma_t}{2}) \|\nabla f(x_t)\|^2 \geq \|\nabla f(\overline{x_k})\|^2 \sum_{t=1}^{k} r_t (1 - \frac{Lr_t}{2})$$

todo:

Will be inplemented later or you can pull requests my [Github Repo](Github Repo)

$$\sum d_i = n, d_i >= 1$$