Efficiently    estimate    $\nabla f(\omega)$

, Let   $G(\omega)$   be   an   unbiased   estimator   for   $\nabla f(\omega)$

$$E(G(\omega)) = \nabla f(\omega)$$

$\rightarrow$ Random  Vector.

$$G(\omega) = 2\left(\langle \omega, x_{i^*}\rangle - y_{i^*}\right) x_{i^*}  \quad \text{where} \quad i^* \text{ sampled}$$

at   a   random.

or

for   $k \ll n$

$$G(\omega) = 2 \sum_{i \in S} \left(\langle \omega, x_i\rangle - y_i\right) x_i  \quad \text{where}$$

$$S \subseteq \text{Input Data}$$

$$|S| = k.$$

THEOREM (SGD CONVERGENCE):

$f$ convex and $\beta$-smooth $\quad \eta \leq 1/\beta$

and $\operatorname{var}(G(x)) \leq \sigma^2$

$$\Rightarrow \quad E\left[f(\bar{x}_k)\right] \leq f(x^*) + \frac{\|x_0 - x^*\|_2^2}{2\eta k} + \eta \sigma^2.$$

$\downarrow$

$$\bar{x}_k = \frac{1}{k}(x_1 + \cdots + x_k)$$

error term introduced by using $G$.

---

So ALL THE NEEDED EQUATIONS:

Convexity: $\forall x, y \quad f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$

$\beta$-Smooth: $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2}\|y - x\|_2^2$

SGD: $x_{k+1} = x_k - \eta G(x_k)$

VECTOR: $\|u - v\|^2 = \|u\|^2 + \|v\|^2 - 2\langle u, v \rangle$

$|\langle u, v \rangle| \leq \|u\| \cdot \|v\|$

UNBIASED : $E[G(x)] = \nabla f(x)$
ESTIMATOR

$\Rightarrow E\left[\|G(x)\|_2^2\right] - \|\nabla f(x)\|_2^2 \leq \sigma^2$

$$\sigma^2 \geq E\left[\|G(x) - E[G(x)]\|_2^2\right] = E\left[\|G(x)\|_2^2\right] - \|E[G(x)]\|_2^2$$

$$= E\left[\|G(x)\|_2^2\right] - \|\nabla f(x)\|_2^2$$

## CLAIM 1:

(SGD    function    decrease    inequality    from    smoothness    alone)

$\forall k$

$$E\left[f(x_{k+1})\right] \leq E\left[f(x_k)\right] - \frac{\eta}{2} E\left[\|\nabla f(x_k)\|_2^2\right] + \frac{\eta \sigma^2}{2}$$

From    $\beta$-smoothness

$$f(x_{k+1}) \leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle +$$

$$\frac{\beta}{2} \|x_{k+1} - x_k\|_2^2$$

$$= f(x_k) + \langle \nabla f(x_k), -\eta G(x_k) \rangle$$

$$+ \frac{\beta \eta^2}{2} \|G(x_k)\|_2^2 \quad - ①$$

We    also    know

$$E\left[\|G(x)\|_2^2\right] - \|\nabla f(x)\|_2^2 \leq \sigma^2$$

From ① , apply expectation $E[G(x_k)] = E[\nabla f(x_k)]$

$$E[f(x_{k+1})] \leq E[f(x_k)] - \eta E[\|\nabla f(x_k)\|_2^2]$$

$$+ \frac{\beta \eta^2}{2}[\sigma^2 + E[\|\nabla f(x)\|_2^2]]$$

$$E[f(x_{k+1})] \leq E[f(x_k)] - \eta E[\|\nabla f(x_k)\|_2^2]$$

$$+ \frac{\beta \eta^2}{2} E[\|\nabla f(x_k)\|_2^2] + \frac{\beta \eta^2 \sigma^2}{2}$$

Given $\beta \leq 1/\eta$

$$\boxed{E[f(x_{k+1})] \leq E[f(x_k)] - \frac{\eta}{2} E[\|\nabla f(x_k)\|_2^2] + \frac{\eta}{2}\sigma^2} \quad - ②$$

Now LET us USE Convexity:

$$f(x_k) \leq f(x^*) + \langle \nabla f(x_k), x_k - x^* \rangle$$

Before that,

$$\|x_{k+1} - x^*\|_2^2 = \|x_k - x^* - \eta G(x_k)\|_2^2 = \|x_k - x^*\|_2^2$$

$$+ \eta^2 \|G(x_k)\|_2^2 - 2\eta \langle x_k - x^*, G(x_k) \rangle$$

$$\left( \|x_{k+1} - x^*\|_2^2 - \|x_k - x^*\|_2^2 \right) = \eta^2 \|G(x_k)\|_2^2 - 2\eta \langle G(x)_k,$$

$$x_k - x^* \rangle$$

Apply expectation

$$E\left( \|x_{k+1} - x^*\|_2^2 - \|x_k - x^*\|_2^2 \right) = \eta^2 E\left( \|G(x_k)\|_2^2 \right) -$$

$$2\eta E \langle G(x)_k, x_k - x^* \rangle$$

$$E\left[ \|G(x)\|_2^2 \right] - E\left[ \|\nabla f(x)_k\|_2^2 \right] \leq \sigma^2 \quad \text{and} \quad E(G(x)) = E(\nabla f(x)_k)$$

$$E\left( \|x_{k+1} - x^*\|_2^2 - \|x_k - x^*\|_2^2 \right) = \eta^2 \sigma^2 + \eta^2 E\left[ \|\nabla f(x)_k\|_2^2 \right]$$

$$2\eta E \left( \langle \nabla f(x)_k, x_k - x^* \rangle \right)$$

Apply  convexity  now

$$f(x_k) \leq f(x^*) + \langle \nabla f(x_k), x_k - x^* \rangle$$

$$-\langle \nabla f(x_k), x_k - x^* \rangle \leq f(x^*) - f(x_k)$$

$$-\mathbb{E}\langle \nabla f(x_k), x_k - x^* \rangle \leq f(x^*) - \mathbb{E}[f(x_k)]$$

$$\mathbb{E}\left[\|x_{k+1} - x^*\|_2^2 - \|x_k - x^*\|_2^2\right] \leq \eta^2 \sigma^2 + \eta^2 \mathbb{E}\left[\|\nabla f(x_k)\|_2^2\right]$$

$$+ \left[f(x^*) - \mathbb{E}[f(x_k)]\right] 2\eta$$

$$\frac{1}{2\eta} \mathbb{E}\left[\|x_{k+1} - x^*\|_2^2 - \|x_k - x^*\|_2^2\right] \leq \frac{\eta}{2}\sigma^2 + \frac{\eta}{2}\mathbb{E}\left[\|\nabla f(x_k)\|_2^2\right]$$

$$+ f(x^*) - \mathbb{E}[f(x_k)]$$

From ②

$$\mathbb{E}[f(x_{k+1})] \leq \mathbb{E}[f(x_k)] - \frac{\eta}{2}\mathbb{E}[\|\nabla f(x_k)\|_2^2] + \frac{\eta}{2}\sigma^2$$

) NO!

$$\frac{1}{2\eta} \mathbb{E}\left[\|x_{k+1} - x^*\|_2^2 - \|x_k - x^*\|_2^2\right] \leq f(x^*) - \mathbb{E}[f(x_{k+1})] + \eta\sigma^2$$

$$\mathbb{E}[f(x_i)] \leq f(x^*) - \frac{1}{2\eta} \mathbb{E}\left[\| x_i - w^* \|_2^2 - \| x_{i-1} - x^* \|_2^2\right]$$

$$+ \eta \sigma^2.$$

**Sum Them:**

$$\sum_{i=1}^{k} \mathbb{E}[f(x_i)] \leq k f(x^*) - \frac{1}{2\eta}\left(\| x_k - x^* \|_2^2 - \| x_0 - x^* \|_2^2\right)$$

$$+ k \eta \sigma^2$$

$$\leq k f(x^*) + \frac{\| x_0 - x^* \|_2^2}{2\eta} + k \eta \sigma^2$$

$$\mathbb{E}[f(\bar{x}_k)] \leq f(x^*) + \frac{\| x_0 - x^* \|_2^2}{2\eta k} + \eta \sigma^2. \quad -\circled{3}$$

Choose $\eta$ such that

$$\frac{\| x_0 - x^* \|_2^2}{2\eta k} = \eta \sigma^2$$

$$\eta = \frac{\| x_0 - x^* \|}{\sigma \sqrt{2k}}$$

Using this $\eta$ in ③

$$E\left[f(\bar{x}_k)\right] \leq f(x^*) + \frac{\eta \|x_0 - x^*\| \sigma \sqrt{2}}{\sqrt{k}}$$

So $\quad \eta \propto \frac{1}{\sqrt{k}}$ .

# Constrained Optimization :

Sometime    other    than

$$\arg\min_x L(x)$$

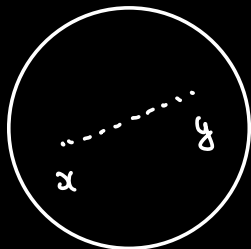we    need    $x$    to    lie    in    some

bounded    region    of    $\mathbb{R}^d$.

$\Rightarrow$ We    can    adapt    GD/SGD) NAND    to    such    situations

as    long    as    constrained    region    is    a    "convex set".

Definition :    $C \subseteq \mathbb{R}^d$    is    a    convex    set    if

$$\forall x, y \in C \ , \ \frac{x+y}{2} \in C \quad \text{(mid point also in c)}$$

equivalently    $\forall x, y \in C \ , \ \forall \lambda \in [0,1] \ , \ \lambda x + (1-\lambda) y \in C.$



C, convex set



not convex set.

## PROJECTED GRADIENT DESCENT:

GOAL:

Compute $\text{argmin}_{x \in C} L(x)$

PROJECTION: $\text{proj}_C(y) = \text{argmin}_{x \in C} \|x - y\|_2$

(closest point in $C$ to $y$)

PGD: $x_{k+1} = \text{Proj}_C(x_k - \eta \nabla f(x_k))$

So, apply gradient descent. If new point $x_{k+1}$ not in $C$, project it into $C$.