Randomized Coordinate Descent $\quad\quad \min\limits_{x \in \mathbb{R}^n} f(x)$

$$\nabla f(x)^{(i)} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \nabla f(x)_i \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

① Gradient Descent is $\quad x_{k+1} = x_k - \eta \nabla f(x_k)$

② Coordinate Descent $\quad x_{k+1} = x_k - \eta \nabla f(x_k)^{i(k)}$

$\quad i(1) = 1$

$\quad i(2) = 2$ $\quad\quad$ only $i(k)$-th entry of $x_k$ is updated

$\quad \vdots$

③ Randomized Coordinate Descent

$$x_{k+1} = x_k - \eta \nabla f(x_k)^{i(k)}$$

$\quad i(k) \sim i.i.d.$ uniform distribution in $\{1, 2, \cdots, n\}$

$\quad\quad$ identical
$\quad\quad$ independent $\quad\quad$ $Prob(i(1) = 1) = \frac{1}{n}$
$\quad\quad$ distributed

Consider an operator $T: \mathbb{R}^n \longrightarrow \mathbb{R}^n$

$$x \longmapsto T(x)$$

and a fixed point iteration

$$x_{k+1} = T(x_k)$$

① An operator $S$ is nonexpansive if

$$\| S(x) - S(y) \| \leq \| x - y \|$$

Example : If $\nabla f(x)$ is $L$-cont. and $f(x)$ is convex

then $S(x) = \left[I - \frac{2}{L}\nabla f\right](x)$ is nonexpansive

$\|S(x) - S(y)\|^2 = \left\| -\frac{2}{L}\left[\nabla f(x) - \nabla f(y)\right] + (x-y)\right\|^2$

$= \|x-y\|^2 + \frac{4}{L^2}\|\nabla f(x) - \nabla f(x)\|^2 - \frac{4}{L}\langle x-y, \nabla f(x) - \nabla f(y)\rangle$

$\leq \|x-y\|^2$

## 2.2.3 Convergence for convex functions

**Theorem 2.8.** *Assume $\nabla f(\mathbf{x})$ is Lipschitz-continuous with Lipschitz constant $L$ and $f(\mathbf{x}) : \mathbb{R}^n \longrightarrow \mathbb{R}$ is convex. Then for any $\mathbf{x}, \mathbf{y}$:*

*1.* $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x}\rangle + \frac{1}{2L}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2$

*2.* $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \leq L\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y}\rangle.$

② $\quad T = (1-\theta)I + \theta S \quad$ with $\theta \in (0,1)$

is called $\theta$-averaged if $S$ is nonexpansive

Example : $S = \left[I - \frac{2}{L}\nabla f\right]$

$T = I - \eta \nabla f = (1-\theta)I + \theta\left(I - \frac{2}{L}\nabla f\right)$

$\theta = \frac{\eta L}{2} \in (0,1) \iff 0 < \eta < \frac{2}{L}$

③ Recall we did the following on Mar 3 :

Theorem (Browder - Göhde - Kirk)

$S : \mathbb{R}^n \to \mathbb{R}^n$ is nonexpansive $\Rightarrow$ $S$ has at least one fixed point.

$S(x_*) = X_*.$

$x_{k+1} = S(x_k)$ may not converge to $x_*$

Example: $S(x) = -x$    $x_* = 0$

**Theorem** If $S: \mathbb{R}^n \to \mathbb{R}^n$ is nonexpansive, then

$$x_{k+1} = (1-\theta)x_k + \theta S(x_k), \quad 0 < \theta < 1$$

converges to one fixed point of $S(x)$.

Example: This implies GD converges if $\eta < \frac{2}{L}$.

④ $\quad T(x) = \begin{bmatrix} [T(x)]_1 \\ [T(x)]_2 \\ \vdots \\ [T(x)]_n \end{bmatrix} \qquad T_i(x) = \begin{bmatrix} x_1 \\ \vdots \\ x_{i-1} \\ [T(x)]_i \\ x_{i+1} \\ \vdots \\ x_n \end{bmatrix}$

If $\quad x_{k+1} = x_k - \eta \nabla f(x_k) = T(x_k)$ is GD

$$x_{k+1} = x_k - \eta \nabla f(x)^{i(k)} \iff x_{k+1} = T_{i(k)}(x_k)$$

**Theorem** Assume

① $T$ is $\theta$-averaged ($\Rightarrow T$ has at least one fixed point)

② $i(k) \in \{1, \cdots, n\}$ is i.i.d. with uniform probability.

Then $\quad x_{k+1} = T_{i(k)}(x_k)$

converges to one fixed point of $T(x)$
with probability 1.

Example: $S = I - \frac{2}{L}\nabla f$ is nonexpansive if $f(x)$ is convex

$$\Rightarrow T = I - \eta \nabla f$$
$$= (1-\theta)I + \theta S \text{ is } \theta\text{-averaged}$$
$$\text{if } \theta = \frac{\eta L}{2} < 1 \iff \eta < \frac{2}{L}$$

Proof: Define $R$, $R_i$ by $\begin{cases} T = I - \theta R \\ T_i = I - \theta R_i \end{cases}$

$$\Rightarrow R_i(x) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ [R(x)]_i \\ 0 \\ \vdots \\ 0 \end{bmatrix} \qquad R = \frac{\eta}{\theta}\nabla f = \frac{2}{L}\nabla f$$
$$T = I - \theta \cdot \frac{2}{L}\nabla f$$

$$x_{k+1} = T_{i(k)}(x_k) \iff x_{k+1} = x_k - \theta R_{i(k)}[x_k]$$

$T$ is $\theta$-averaged $\iff T = (1-\theta)I + \theta S$
with $S$ nonexpansive

$I - R = I - \frac{I-T}{\theta}$
$\iff \frac{1}{\theta}T - (\frac{1}{\theta}-1)I$ is nonexpansive

$= (1-\frac{1}{\theta})I + \frac{1}{\theta}T$
$\iff I - R$ is nonexpansive

$= S$
$\iff \|x - Rx - y + Ry\|^2 \leq \|x-y\|^2$

$\iff \frac{1}{2}\|Rx - Ry\|^2 \leq \langle x-y, Rx-Ry\rangle$

$T = I - \theta R$

$T(x_*) = x_* \iff R(x_*) = 0$

$y = x_* \Rightarrow \frac{1}{2}\|Rx\|^2 \leq \langle Rx, x - x_*\rangle$

$\frac{1}{2}\|\frac{2}{L}\nabla f(x)\|^2 \leq \langle \frac{2}{L}\nabla f(x), x - x_*\rangle$

$\iff \|\nabla f(x)\|^2 \leq L\langle \nabla f(x), x - x_*\rangle$

Definition of Expectation & more:

Example: ① $X$ is a random variable taking values in $\{0,1\}$

with equal probability

definition $\begin{cases} E(X) = 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{1}{2} \\ E(X^2) = 0^2 \cdot \frac{1}{2} + 1^2 \cdot \frac{1}{2} = \frac{1}{2} \\ E(f(X)) = f(0) \cdot \frac{1}{2} + f(1) \cdot \frac{1}{2} = \frac{1}{2} f(1) \end{cases}$

$P(X=0) = \frac{1}{2}$

$P(X=1) = \frac{1}{2}$

$E$ denotes expectation w.r.t. random variable $X$

Expectation for discrete random variables are defined as above

② $X$ is a random variable taking values

in $\{X_1, X_2, \cdots, X_k\}$

with probability $P_1, P_2, \cdots, P_k$     $\sum_{i=1}^{k} P_i = 1$, $P_i \geq 0$

definition $\begin{cases} E(X) = X_1 \cdot P_1 + X_2 \cdot P_2 + \cdots + X_k P_k \rightarrow \text{Convex} \\ E(f(X)) = \sum_{i=1}^{k} f(X_i) P_i \end{cases}$ combination

$\text{Prob}(X = X_i) = P_i$

Jensen's: $f(E(X)) \leq E(f(X))$ if $f(x)$ is convex

Example: $f(x) = x^2 \Rightarrow [E(X)]^2 \leq E(X^2)$

③ X, Y are i.i.d. random variable taking values
$$\{a_1, a_2, \ldots, a_k\}$$
with probability $P_1, P_2, \ldots, P_k$

$$\sum_{i=1}^{k} P_i = 1, \quad P_i \geq 0$$

Joint probability $P(X=a_i, Y=a_j)$

X & Y are independent $\Rightarrow$

$$P(X=a_i, Y=a_j) = P(X=a_i) P(Y=a_j)$$

$$E[f(X,Y)] = \sum_{i=1}^{k} \sum_{j=1}^{k} f(a_i, a_j) P(X=a_i, Y=a_j)$$

$$= \sum_{i=1}^{k} \sum_{j=1}^{k} f(a_i, a_j) P_i P_j$$

In $\quad X_{k+1} = T_{i(k)}(X_k)$

$$X_{k+1} = X_k - \eta \nabla f(X_k)^{i(k)}$$

$X_0$ is deterministic and $i(0), i(1), \ldots, i(N)$ are random

$X_N$ is a function of these $N+1$ random variables

$E(X_N)$ denotes expectation w.r.t. $N+1$ i.i.d. random variables

④ Conditional Probability & Expectation

X is a random variable taking values

in $\{X_1, X_2, \cdots, X_k\}$

with probability $P_1, P_2, \cdots, P_k$

$Y$ is a random variable taking values

in $\{y_1, y_2, \cdots, y_\ell\}$

with probability $q_1, q_2, \cdots, q_\ell$

Probability of event $X = X_i$ given the knowledge $Y = y_j$

$$P(X = X_i \mid Y = y_j) = \frac{P(X = X_i, Y = y_j)}{P(Y = y_j)} \leftarrow \text{joint prob.}$$

The conditional expectation w.r.t. $X$ given $Y = y_j$

$$E(X \mid Y = y_j) = \sum_{i=1}^{k} X_i \, P(X = X_i \mid Y = y_j)$$

$$= \sum_{i=1}^{k} X_i \, \frac{P(X = X_i, Y = y_j)}{P(Y = y_j)}$$

$$g(y) = E(X \mid Y = y) = \sum_{i=1}^{k} X_i \, P(X = X_i \mid Y = y)$$

is $\begin{cases} \text{a function of } y \\ \text{a random variable} \end{cases}$

We can also write it as $g(Y) = E(X \mid Y)$

So $E(X \mid Y)$ is a random variable.

$$E_k = E\left[ i(k) \mid i(k-1), i(k-2), \cdots, i(0) \right]$$

$E_k$ denotes the conditional expectation w.r.t. $i(k)$

Conditioned on the past random variables

$$\hat{i}(k-1), \hat{i}(k-2), \dots, \hat{i}(0)$$

$$X_{k+1} = T_{\hat{i}(k)}(X_k)$$
$$X_{k+1} = X_k - \eta \nabla f(X_k)^{\hat{i}(k)}$$

Then ① $\underbrace{E_k(X_k)}_{\downarrow} = X_k$ because $X_k$ does

random variable $\quad$ NOT depend on $\hat{i}(k)$

② Let $X$ be something depending on all $\hat{i}(k)$
$$\hat{i}(k-1), \hat{i}(k-2), \dots, \hat{i}(0)$$

$$E[E_k(X)] = E[X]$$

Law of Total Expectation : $E[E(X|Y)] = E(X)$

$E(X|Y)$ is a function of $Y$ (random variable)

$$E(X|Y) = \sum_{i=1}^{k} x_i P(X = x_i | Y)$$

$$E[E(X|Y)] = E\left[\sum_{i=1}^{k} x_i P(X = x_i | Y)\right]$$

$$= \sum_{j=1}^{\ell} \left[\sum_{i=1}^{k} x_i P(X = x_i | Y = y_j)\right] \cdot P(Y = y_j)$$

$$= \sum_{j=1}^{\ell} \sum_{i=1}^{k} x_i P(X = x_i | Y = y_j) \cdot P(Y = y_j)$$

$$= \sum_{j=1}^{\ell} \sum_{i=1}^{k} x_i P(X = x_i, Y = y_j)$$

$$= \sum_{i=1}^{k} x_i \left[\sum_{j=1}^{\ell} P(X = x_i, Y = y_j)\right]$$

$$= \sum_{i=1}^{k} x_i \, P(x = x_i)$$

$$= E(x)$$

③ 
$$E_k \left[ R_{\hat{i}(k)} (x_k) \right] = \frac{1}{n} R(x_k)$$

$$R = \frac{2}{L} \nabla f \qquad R_i(x) = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \nabla f(x)_i \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$E_k$ is expectation w.r.t. $\hat{i}(k)$

$$\begin{cases} P(\hat{i}(k) = 1) = \frac{1}{n} \\ P(\hat{i}(k) = 2) = \frac{1}{n} \\ \quad \vdots \\ P(\hat{i}(k) = n) = \frac{1}{n} \end{cases}$$

$$\Rightarrow E_k \left[ R_{\hat{i}(k)} (x_k) \right] = \sum_{j=1}^{n} P[\hat{i}(k) = j] \cdot R_j(x_k)$$

$$= \frac{1}{n} \sum_{j=1}^{n} R_j(x_k) = \frac{1}{n} R(x_k)$$

$$E_k \, \| R_{\hat{i}(k)} (x_k) \|^2 = \frac{1}{n} \| R(x_k) \|^2$$

$$\| R_{i(k)} (x_k) \|^2 = \left( [R(x_k)]_{\hat{i}(k)} \right)^2$$

$$E_k \left( [R(x_k)]_{\hat{i}(k)} \right)^2 = \sum_{j=1}^{n} P[\hat{i}(k) = j] \cdot \left( R_j(x_k) \right)^2$$

$$= \frac{1}{n} \sum_{j=1}^{n} \left( R_j (x_k) \right)^2 = \frac{1}{n} \| R(x_k) \|^2$$

④ Let $x_*$ be a fixed point to $x_* = T(x_*)$

$$\| x_{k+1} - x_* \|^2 = \| x_k - \theta R_{i(k)} (x_k) - x_* \|^2$$

$$= \| x_k - x_* \|^2 - 2\theta \langle R_{i(k)} x_k, x_k - x_* \rangle$$
$$+ \theta^2 \| R_{i(k)} (x_k) \|^2$$

Take $E_k$ ( conditional expectation w.r.t. $i(k)$ )

Expectation is linear

$$E_k \| x_{k+1} - x_* \|^2 = E_k \| x_k - x_* \|^2 - 2\theta E_k \langle R_{i(k)} x_k, x_k - x_* \rangle$$
$$+ \theta^2 E_k \| R_{i(k)} (x_k) \|^2$$

$$= \| x_k - x_* \|^2 - 2\theta \langle E_k [ R_{i(k)} (x_k) ], x_k - x_* \rangle$$
$$+ \theta^2 \cdot \frac{1}{n} \| R(x_k) \|^2$$

$$= \| x_k - x_* \|^2 - \frac{2\theta}{n} \langle R(x_k), x_k - x_* \rangle$$
$$+ \theta^2 \cdot \frac{1}{n} \| R(x_k) \|^2$$

$$\leq \| x_k - x_* \|^2 - \frac{\theta}{n} \| R(x_k) \|^2 + \frac{\theta^2}{n} \| R(x_k) \|^2$$

$$\boxed{ \frac{1}{2} \| Rx \|^2 \leq \langle Rx, x - x_* \rangle }$$

Take $E$ for both sides

$$\Rightarrow E \| x_{k+1} - x_* \|^2 \leq E \| x_k - x_* \|^2 - (1-\theta) \frac{\theta}{n} E \| R(x_k) \|^2$$
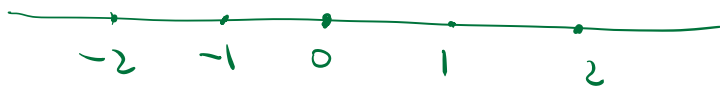
$$\Rightarrow \quad \left\{ E \| X_k - X_* \| \right\} \searrow$$

⑤ Discrete Martingale

is a sequence of random variable $X_1, X_2, X_3, \ldots$ satisfying

① $E(|X_n|) < +\infty$

② $E(X_{n+1} | X_1, \ldots, X_n) = X_n$

Example: positions of 1D random walk

Integer grid; $X_0 = 0$; $X_{k+1} = \begin{cases} X_k + 1 & \text{with prob. } \frac{1}{2} \\ X_k - 1 & \text{with prob. } \frac{1}{2} \end{cases}$

$E(|X_n|) \leq n$

Given $X_1, \ldots, X_n$, expected value of $X_{n+1}$ is $X_n$

$E(X_{n+1} | X_1, \ldots, X_n) = X_n$

Supermartingale if $E(X_{n+1} | X_1, \ldots, X_n) \leq X_n$

submartingale if $E(X_{n+1} | X_1, \ldots, X_n) \geq X_n$

Supermartingale Convergence Theorem

Let $X_k, Y_k$ be random variables satisfying

① $X_k \geq 0, Y_k \geq 0$ almost surely

② $E[X_{k+1} | X_1, \ldots, X_k] \leq X_k - Y_k$

Then almost surely (with probility 1)

1) $X_k$ converges to $X_\infty$, $k \to \infty$

2) $\sum_{k=0}^{\infty} Y_k < +\infty$

Remark: $X_\infty$ is a random variable

$$E_k \|X_{k+1} - X_*\|^2 \leq \|X_k - X_*\|^2 - (1-\theta)\frac{\theta}{n} \|R(X_k)\|^2$$

$\Rightarrow \begin{cases} \sum_{k=0}^{\infty} \|R(X_k)\|^2 < \infty \Rightarrow \|R(X_k)\| \to 0 \\ \lim_{k \to \infty} \|X_k - X_*\|^2 \text{ exists} \end{cases}$  with probability 1

Let $Fix(T)$ be the set of all fixed points of $T$

$\forall X_* \in Fix(T)$, $\left[\lim_{k \to \infty} \|X_k - X_*\| \text{ exists with prob. } 1\right]$

$\Rightarrow$ With prob. 1 $\left[\forall X_* \in Fix(T), \lim_{k \to \infty} \|X_k - X_*\| \text{ exists}\right]$

$\to$ nontrivial, skipped, see last reference book.

⑥ With prob. 1, we have

1) $\{X_k\}$ is a bounded sequence

thus a convergent subsequence $\{X_{k_j}\} \to Z$

2) $\|R(x_k)\| \to 0 \implies \|\frac{1}{\theta}(I-T)(x_k)\| \to 0$

$\implies \|(I-T)x_k\| \to 0$

$\implies \|(I-T)x_{k_j}\| \to 0$

T is $\theta$ averaged

$\implies \|T(x) - T(y)\| \leqslant \|x-y\|$

$\implies$ T is continuous

$\implies I-T$ is continuous

$\implies \|z - T(z)\| = 0$

$\implies z - T(z) = 0$

$\implies z \in Fix(T)$

$$\|x_{k+1} - x_*\|^2 = \|x_k - x_*\|^2 - 2\theta \langle R_{i(k)} x_k, x_k - x_* \rangle + \theta^2 \|R_{i(k)}(x_k)\|^2$$

$\boxed{\|R(x_k)\| \to 0}$

set $x_* = z \to \{x_k\}$ has the same limit as $x_{k_j}$