

Digest of the convergence analysis of APG by Xu-Yin

Problem (1.1)

$$\min_{x \in \mathcal{X}} F(x_1, \dots, x_S) \triangleq f(x_1, \dots, x_S) + \sum_{i=1}^S r_i(x_i) \quad (1.1)$$

- x_i block variable, x full variable, S -blocks
- \mathcal{X} constraint on x , assumed:
 - close
 - block convex ... (def 1)
- f assumed block convex, differentiable, smooth
- r_i extended valued $(r_i(x_i) = \infty \text{ if } x_i \notin \text{dom } r_i)$
 - r_i can be non-smooth

Note: r_i, f can be non-convex

Interpretation of (1.1)

- f models the joint function, r_i model individual block
- r_i models individual constraint
- \mathcal{X} models joint constraint
- if no constraint on block i , $r_i \equiv 0$

Definitions

(def 1) Set multi-convex (mc) = the set \mathcal{X} is mc if its projection to each sub-block is cvx.

i.e. $\forall i$, fixed $S-1$ block $\{x_j\}_{j \neq i}$, the set

$$\mathcal{X}_i(\{x_j\}_{j \neq i}) \triangleq \left\{ x_i \in \mathbb{R}^{n_i} \mid \underbrace{(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_S)}_{\text{fix}} \in \mathcal{X} \right\}$$

↑
subject

is cvx.

(def2) Block multi-convex function = $\forall i$, f is convex of x_i while other block held fix.

i.e. for C.D, the function $F(x_1, \dots, x_s)$ is convex w.r.t. one block while other block are fix.

BCD update of Gauss-Seidel Type / Cyclic update

"minimizing F cyclically over x_1, x_2, \dots, x_s while fixing other blocks at their last updated values"

Notation $x_i^k = x_i$ at k^{th} iteration

$$f_i^k(x_i) \triangleq f(\underbrace{x_1^k, \dots, x_{i-1}^k}_{\text{fix, at last updated value}}, \overset{\text{subject}}{x_i}, \underbrace{x_{i+1}^{k-1}, \dots, x_s^{k-1}}_{\text{fix, at last updated value}})$$

note = k and $k-1$!!

$$g_i^k = \nabla f_i^k(x_i^{k-1}) \quad \text{block partial gradient of } f \text{ at } x_i^{k-1}$$

note: g_i^k and x_i^{k-1} !!

$$\hat{x}_i^{k-1} = x_i^{k-1} + \omega_i^{k-1} (x_i^{k-1} - x_i^{k-2})$$

extrapolation point, $\omega_i^{k-1} \geq 0 \forall i, k$, the weight

Update forms

(a) original update
$$x_i^k = \operatorname{argmin}_{x_i \in X_i^k} f_i^k(x_i) + r_i(x_i)$$

"update x_i based on $f_i^k(x_i)$ and r_i "

(b) proximal update
$$x_i^k = \operatorname{argmin}_{x_i \in X_i^k} f_i^k(x_i) + \frac{L_i^{k-1}}{2} \|x_i - x_i^{k-1}\|^2 + r_i(x_i)$$

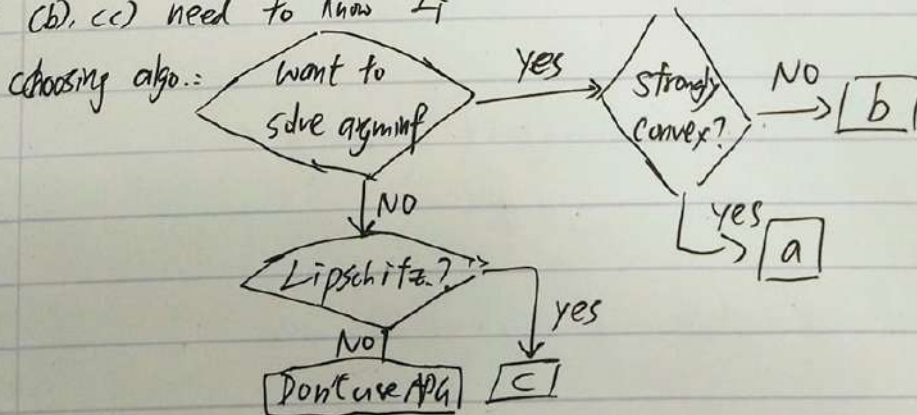
"update x_i based on original form plus a proximal term, where block-Lipschitz constant L_i^{k-1} is known"

(c) prox-linear
$$x_i^k = \operatorname{argmin}_{x_i \in X_i^k} \{ \hat{g}_i^k(x_i^{k-1}), x_i - x_i^{k-1} \} + \frac{L_i}{2} \|x_i - x_i^{k-1}\|^2 + r_i(x_i)$$

"replace f in (b) by linear approximation"

Remarks

- if f, X convex, then all sub-problems of (a), (b), (c) are convex \Rightarrow global minima for all sub-problem
- in general harder to solve (a), (b) than (c) due to the "argmin f "
- (b), (c) need to know L_i^{k-1}



Assumptions

(a-1) F is continuous in dom F (a-2) $\inf_{x \in \text{dom} F} F(x) > -\infty$

(a-3) (1.1) has a Nash-point (useless for matrix factorization problem)

(a-4) Consistent update scheme (a,b,c) for each block

(a-5) For (a), f_i^k is str-conv: $l_i \leq L_i^{k-1} \leq L_i$ (a-6) For (b), L_i^{k+1} obey $l_i \leq L_i^{k+1} \leq L_i$ (a-7) For (c), $\forall f_i^k$ is Lipschitz L_i^{k+1} obey $l_i \leq L_i^{k+1} \leq L_i$

$$f_i^k(x_i^k) \leq f_i^k(x_i^{k+1}) + \langle \bar{g}_i^k, x_i^k - x_i^{k+1} \rangle + \frac{L_i^{k-1}}{2} \|x_i^k - x_i^{k+1}\|_2^2$$

Conceptual flow of the convergence proof

Step I 1. Show sufficient decrease condition

2. Show square summable

$$\sum_k \|x^k - x^{k+1}\|_2^2 < +\infty$$

Step 2 1. Use Kurdyka-Łojasiewicz inequality to improve $\sum \| \cdot \|_2^2 \rightarrow \sum \| \cdot \|_2^1$

2. Global convergence + asymptotic rate

Lemma 2.1. Let $a(u), b(u)$ be convex function on convex set U .

a is differentiable, b is possibly non-smooth. ... (A)
Let $A(u) = a(u) + b(u)$.

Let $u^* = \operatorname{argmin}_{u \in U} \langle a(v), u-v \rangle + \frac{1}{2} \|u-v\|^2 + b(u)$ (Δ)

If $a(u^*) \leq a(v) + \langle \nabla a(v), u^*-v \rangle + \frac{1}{2} \|u^*-v\|^2$... (#)

Then

$$A(u) - A(u^*) \geq \frac{L}{2} \|u^*-v\|^2 + L \langle v-u^*, u^*-v \rangle \quad \forall u$$

Proof. From (Δ), by first-order optimality / Fermat's rule on u^* ,

$$\langle \nabla a(v) + L(u^*-v) + g, u-u^* \rangle \geq 0, \quad g \in \partial b(u^*) \quad \dots (*)$$

(why = recall 1. b is possibly non-smooth, so use ∂b
2. $u^* = \operatorname{argmin} \phi \Leftrightarrow \langle \nabla \phi(u^*), u-u^* \rangle \geq 0 \quad \forall u$)

Consider $A(u) - A(u^*)$

$$A(u) - A(u^*) \stackrel{(A)}{=} a(u) - a(u^*) + b(u) - b(u^*)$$

$$\stackrel{(\#)}{\geq} a(u) - \left(a(v) + \langle \nabla a(v), u^*-v \rangle + \frac{L}{2} \|u^*-v\|^2 \right) + b(u) - b(u^*)$$

$$= a(u) - a(v) - \left(\langle \nabla a(v), u^*-v \rangle + \frac{L}{2} \|u^*-v\|^2 \right) + b(u) - b(u^*)$$

Trick $-\langle \nabla a(v), u^*-v \rangle = -\langle \nabla a(v), u^*-v + \underline{u-u} \rangle$
 $= -\langle \nabla a(v), u^*-u \rangle - \langle \nabla a(v), u-v \rangle$
 $= + \langle \nabla a(v), u-u^* \rangle - \langle \nabla a(v), u-v \rangle$

Trick $\Rightarrow \underline{a(u) - a(v)}$

$$\begin{aligned} & + \langle a(v), u-u^* \rangle - \langle \nabla a(v), u-v \rangle - \frac{L}{2} \|u^*-v\|^2 \\ & + b(u) - b(u^*) \\ \stackrel{a \text{ is convex}}{\geq} & \langle a(v), u-u^* \rangle - \frac{L}{2} \|u^*-v\|^2 + b(u) - b(u^*) \end{aligned}$$

6

Now we have

$$A(u) - A(u^*) \geq \langle \nabla a(v), u - u^* \rangle - \frac{L}{2} \|u^* - v\|_2^2 + b(u) - b(u^*)$$

To handle $\langle \nabla a(v), u - u^* \rangle$, use (*)

$$(*) \Rightarrow \langle \nabla a(v), u - u^* \rangle \geq -L \langle u^* - v, u - u^* \rangle - \langle g, u - u^* \rangle$$

$$\text{So } A(u) - A(u^*) \stackrel{(*)}{\geq} -L \langle u^* - v, u - u^* \rangle - \langle g, u - u^* \rangle - \frac{L}{2} \|u^* - v\|_2^2 + \underbrace{b(u) - b(u^*)}$$

*) b convex

$$\geq -L \langle u^* - v, u - u^* \rangle - \frac{L}{2} \|u^* - v\|_2^2$$

$$= L \langle u^* - v, u^* - u + v - v \rangle - \frac{L}{2} \|u^* - v\|_2^2$$

$$= \frac{L}{2} \|u^* - v\|_2^2 + L \langle u^* - v, v - u \rangle //$$

Lemma 2.2

1. Assumption 1.2
2. $\{x^k\}$ generate by algo-1
3. With $0 \leq \omega_i^{k-1} \leq \delta \omega \sqrt{\frac{L_i^{k-2}}{L_i^{k-1}}}$ $\forall i \in I_3$, $\delta \omega < 1$

$$\text{Then } \sum_{k=0}^{\infty} \|x^k - x^{k+1}\|_2^2 < \infty$$

proof. For I_3 , let $a = f_i^k$, $b = r_i$, $v = \bar{x}_i^{k-1}$, $u = x_i^{k-1}$ ($\Rightarrow u^* = x_i^k$)
 $F = f_i + r_i$

$$F_i^k(x_i^{k-1}) - F_i^k(x_i^k)$$

lemma 2.1

$$\geq \frac{L_i^{k-1}}{2} \|x_i^k - \bar{x}_i^{k-1}\|_2^2 + L_i^{k-1} \langle \bar{x}_i^{k-1} - x_i^{k-1}, x_i^k - \bar{x}_i^{k-1} \rangle$$

$$F_i^K(x_i^{K+1}) - F_i^K(x_i^K) \geq \underbrace{\frac{L_i^{K-1}}{2} \|x_i^K - \tilde{x}_i^{K-1}\|_2^2}_{\textcircled{1}} + L_i^{K-1} \underbrace{\langle \tilde{x}_i^{K-1} - x_i^{K-1}, x_i^K - \tilde{x}_i^{K-1} \rangle}_{\textcircled{2}}$$

$$\begin{aligned} \textcircled{1} &= \|x_i^K - \tilde{x}_i^{K-1} + \tilde{x}_i^{K-1} - x_i^{K-1}\|_2^2 \\ &= \|x_i^K - \tilde{x}_i^{K-1}\|_2^2 + \|\tilde{x}_i^{K-1} - x_i^{K-1}\|_2^2 + 2\langle x_i^K - \tilde{x}_i^{K-1}, \tilde{x}_i^{K-1} - x_i^{K-1} \rangle \end{aligned}$$

(how to think: expand it to get terms to cancel with $\textcircled{2}$)

$$\begin{aligned} \frac{L_i^{K-1}}{2} \textcircled{1} + L_i^{K-1} \textcircled{2} &= \frac{L_i^{K-1}}{2} \|x_i^K - \tilde{x}_i^{K-1}\|_2^2 \\ &\quad + \frac{L_i^{K-1}}{2} \|\tilde{x}_i^{K-1} - x_i^{K-1}\|_2^2 \\ &\quad - L_i^{K-1} \|\tilde{x}_i^{K-1} - x_i^{K-1}\|_2^2 \end{aligned} \left. \vphantom{\frac{L_i^{K-1}}{2} \textcircled{1} + L_i^{K-1} \textcircled{2}} \right\} -\frac{L_i^{K-1}}{2} \|\tilde{x}_i^{K-1} - x_i^{K-1}\|_2^2$$

$$\text{Hence: } F_i^K(x_i^{K+1}) - F_i^K(x_i^K) \geq \frac{L_i^{K-1}}{2} \|x_i^K - \tilde{x}_i^{K-1}\|_2^2 - \frac{L_i^{K-1}}{2} \|\tilde{x}_i^{K-1} - x_i^{K-1}\|_2^2$$

Now cancel \tilde{x}_i^{K-1} via update $\tilde{x}_i^{K-1} = x_i^{K-1} + \omega_i^{K-1} (x_i^{K-1} - x_i^{K-2})$

$$\begin{aligned} F_i^K(x_i^{K+1}) - F_i^K(x_i^K) &\geq \frac{L_i^{K-1}}{2} \|x_i^K - x_i^{K-1}\|_2^2 \\ &\quad - \frac{L_i^{K-1}}{2} (\omega_i^{K-1})^2 \|x_i^{K-1} - x_i^{K-2}\|_2^2 \end{aligned}$$

$$\begin{aligned} &\stackrel{\textcircled{a}}{\geq} \frac{L_i^{K-1}}{2} \|x_i^K - x_i^{K-1}\|_2^2 \\ &\quad - \frac{L_i^{K-2}}{2} \rho_w^2 \|x_i^{K-1} - x_i^{K-2}\|_2^2 \quad \dots (2.8) \end{aligned}$$

$$\textcircled{a}) \omega_i^{K-1} \leq \rho_w \sqrt{\frac{L_i^{K-2}}{L_i^{K-1}}} \Leftrightarrow -(\omega_i^{K-1})^2 \leq -\rho_w^2 \frac{L_i^{K-2}}{L_i^{K-1}}$$

(2.8) holds for update (c), but also (a), (b) as

$$F_i^K(x_i^{K-1}) - F_i^K(x_i^K) \geq \frac{L_i^{K-1}}{2} \|x_i^{K-1} - x_i^K\|_2^2$$

"Sufficient decrease condition"!!

(2.8) is on block i , now extend to all block

$$F(x^{k+1}) - F(x^k) = \sum_{i=1}^S \left(F_i^k(x_i^{k+1}) - F_i^k(x_i^k) \right)$$

$$\stackrel{(2.8)}{\geq} \sum_{i=1}^S \left(\frac{L_i^{k+1}}{2} \|x_i^{k+1} - x_i^k\|_2^2 - \frac{L_i^{k+1} \rho_{\omega}^2}{2} \|x_i^{k+1} - x_i^{k-1}\|_2^2 \right)$$

Telescoping sum:

$$\begin{aligned} F(x^0) - F(x^1) &= \text{---} \\ F(x^1) - F(x^2) &= \text{---} \\ &\vdots \\ \textcircled{\oplus} F(x^{k+1}) - F(x^k) &= \text{---} \end{aligned}$$

$$F(x^0) - F(x^k) = \sum \sum (\cdot)$$

$$F(x^0) - F(x^k) \geq \sum_{k=1}^K \sum_{i=1}^S \left(\frac{L_i^{k+1}}{2} \|x_i^{k+1} - x_i^k\|_2^2 - \frac{L_i^{k+1} \rho_{\omega}^2}{2} \|x_i^{k+1} - x_i^{k-1}\|_2^2 \right)$$

expand on k

$$= \sum_{i=1}^S \left(\frac{L_i^0}{2} \|x_i^0 - x_i^1\|_2^2 - \frac{L_i^1 \rho_{\omega}^2}{2} \|x_i^1 - x_i^0\|_2^2 \right) \quad 0 \text{ as } x^{k-1} = x^0$$

$$+ \frac{L_i^1}{2} \|x_i^1 - x_i^2\|_2^2 - \frac{L_i^2 \rho_{\omega}^2}{2} \|x_i^2 - x_i^1\|_2^2$$

$$+ \frac{L_i^2}{2} \|x_i^2 - x_i^3\|_2^2 - \frac{L_i^3 \rho_{\omega}^2}{2} \|x_i^3 - x_i^2\|_2^2$$

+ ...

$$+ \frac{L_i^{k+1}}{2} \|x_i^{k+1} - x_i^k\|_2^2 - \frac{L_i^{k+1} \rho_{\omega}^2}{2} \|x_i^{k+1} - x_i^{k-1}\|_2^2$$

$$= \sum_{i=1}^S \left(\frac{(1 - \rho_{\omega}^2) L_i^0}{2} \|x_i^1 - x_i^0\|_2^2 \right.$$

$$+ \frac{(1 - \rho_{\omega}^2) L_i^1}{2} \|x_i^2 - x_i^1\|_2^2$$

$$+ \frac{(1 - \rho_{\omega}^2) L_i^2}{2} \|x_i^3 - x_i^2\|_2^2$$

$$+ \dots + \frac{(1 - \rho_{\omega}^2) L_i^{k+1}}{2} \|x_i^{k+1} - x_i^k\|_2^2 \Big) + \frac{L_i^{k+1}}{2} \|x^{k+1} - x^k\|_2^2$$

$$\geq \sum_{k=1}^K \sum_{i=1}^S \frac{(1 - \rho_{\omega}^2) L_i^{k+1}}{2} \|x_i^{k+1} - x_i^k\|_2^2$$

$$\geq \sum_{k=1}^K \frac{(1 - \rho_{\omega}^2) l_i}{2} \|x_i^{k+1} - x_i^k\|_2^2 \quad (l_i \leq L_i^{k+1} \forall i)$$

Now we have

$$F(x^0) - F(x^K) \geq \underbrace{\frac{1 - \beta_{\omega}^2}{2}}_c \sum_{i=1}^K \|x^{k-1} - x^k\|^2$$

as F is lower bounded, $-\infty < \inf F \leq F(x^K)$

$$\Rightarrow -F(x^K) < +\infty \Rightarrow F(x^0) - F(x^0) < +\infty$$

$$\Rightarrow +\infty \Rightarrow F(x^0) - F(x^K) \geq c \sum_{k=1}^K \|x^{k-1} - x^k\|^2$$

$$\Rightarrow \sum_{k=1}^K \|x^{k-1} - x^k\|^2 < +\infty$$

$\Rightarrow \|x^{k-1} - x^k\|_2$ is square-summable!

\Rightarrow the sequence $\{x^k\}$ produced by APG converge. //

Note. Lemma 2.2 said $\{x^k\}$ converge.

It does not say where it converge and how fast it converge.

"where" it converge is by the next lemma.

By Kurdyka-Łojasiewicz inequality, we have

$$\sum_{k=1}^K \|x^{k-1} - x^k\|_2^2 < +\infty$$