

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Accelerated Proximal Gradient Methods for Nonconvex Programming

Anonymous Author(s)
Affiliation
Address
email

We consider a general problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x}), \quad (1)$$

We mainly consider nonconvex f and nonconvex nonsmooth g .

1 Preliminaries

1.1 Basic Assumptions

Definition 1 A function $g : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ is said to be proper if $\text{dom } g \neq \emptyset$, where $\text{dom } g = \{\mathbf{x} \in \mathbb{R}^n : g(\mathbf{x}) < +\infty\}$. g is lower semicontinuous at point \mathbf{x}_0 if

$$\liminf_{\mathbf{x} \rightarrow \mathbf{x}_0} g(\mathbf{x}) \geq g(\mathbf{x}_0). \quad (2)$$

In problem (1), we assume that f is a proper function with Lipschitz continuous gradients and g is proper and lower semicontinuous. We assume that $F(\mathbf{x})$ is coercive, i.e., F is bounded from below and

$$F(\mathbf{x}) \rightarrow \infty \quad \text{when} \quad \|\mathbf{x}\| \rightarrow \infty, \quad (3)$$

where $\|\cdot\|$ is the l_2 -norm.

1.2 Subdifferentials of Nonconvex and Nonsmooth Functions

Definition 2 [1, 2] Let g be a proper and lower semicontinuous function.

1. For a given $\mathbf{x} \in \text{dom } g$, the Frechet subdifferential of g at \mathbf{x} , written as $\hat{\partial}g(\mathbf{x})$, is the set of all vectors $\mathbf{u} \in \mathbb{R}^n$ which satisfy

$$\liminf_{\mathbf{y} \neq \mathbf{x}, \mathbf{y} \rightarrow \mathbf{x}} \frac{g(\mathbf{y}) - g(\mathbf{x}) - \langle \mathbf{u}, \mathbf{y} - \mathbf{x} \rangle}{\|\mathbf{y} - \mathbf{x}\|} \geq 0. \quad (4)$$

2. The limiting-subdifferential, or simply the subdifferential, of g at $\mathbf{x} \in \mathbb{R}^n$, written as $\partial g(\mathbf{x})$, is defined through the following closure process

$$\partial f(\mathbf{x}) := \{\mathbf{u} \in \mathbb{R}^n : \exists \mathbf{x}_k \rightarrow \mathbf{x}, g(\mathbf{x}_k) \rightarrow g(\mathbf{x}), \mathbf{u}_k \in \hat{\partial}g(\mathbf{x}_k) \rightarrow \mathbf{u}, k \rightarrow \infty\}. \quad (5)$$

Proposition 1 [1, 2]

1. In the nonsmooth context, the Fermat's rule remains unchanged: If $\mathbf{x} \in \mathbb{R}^n$ is a local minimizer of g , then $0 \in \partial g(\mathbf{x})$.
2. Let $(\mathbf{x}_k, \mathbf{u}_k)$ be a sequence such that $\mathbf{x}_k \rightarrow \mathbf{x}$, $\mathbf{u}_k \rightarrow \mathbf{u}$, $g(\mathbf{x}_k) \rightarrow g(\mathbf{x})$ and $\mathbf{u}_k \in \partial g(\mathbf{x}_k)$, then $\mathbf{u} \in \partial g(\mathbf{x})$.

054 3. If f is a continuously differentiable function, then $\partial(f + g)(\mathbf{x}) = \nabla f(\mathbf{x}) + \partial g(\mathbf{x})$.

055
056 Recall that points whose subdifferential contains 0 are called critical points.

057 058 1.3 Proximal Mapping

059
060 Let $g : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ be a proper and semicontinuous function. Given $\mathbf{x} \in \mathbb{R}^n$ and $\alpha > 0$,
061 define the proximal mapping [1] as:

$$062 \text{prox}_{\alpha g}(\mathbf{x}) = \underset{\mathbf{u}}{\operatorname{argmin}} g(\mathbf{u}) + \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{u}\|^2. \quad (6)$$

063
064
065 When $g := \delta_X$, the indicator function of a nonempty and closed set X , defined as:

$$066 \delta_X(\mathbf{x}) = \begin{cases} 0, & \text{if } \mathbf{x} \in X, \\ \infty, & \text{otherwise,} \end{cases} \quad (7)$$

067
068 the proximal mapping reduces to the projection onto X .

069 070 1.4 KL Inequality

071
072 **Definition 3** [3, 2] A function $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ is said to have the KL property at $\bar{\mathbf{u}} \in$
073 $\operatorname{dom} \partial f := \{\mathbf{x} \in \mathbb{R}^n : \partial f(\mathbf{x}) \neq \emptyset\}$ if there exists $\eta \in (0, +\infty]$, a neighborhood U of $\bar{\mathbf{u}}$ and a
074 function $\varphi \in \Phi_\eta$, such that for all

$$075 \mathbf{u} \in U \cap \{\mathbf{u} \in \mathbb{R}^n : f(\bar{\mathbf{u}}) < f(\mathbf{u}) < f(\bar{\mathbf{u}}) + \eta\}, \quad (8)$$

076
077 the following inequality holds

$$078 \varphi'(f(\mathbf{u}) - f(\bar{\mathbf{u}})) \operatorname{dist}(0, \partial f(\mathbf{u})) > 1, \quad (9)$$

079 where Φ_η stands for a class of function $\varphi : [0, \eta) \rightarrow \mathbb{R}^+$ satisfying: (1) φ is concave and C^1 on
080 $(0, \eta)$; (2) φ is continuous at 0, $\varphi(0) = 0$; and (3) $\varphi'(\mathbf{x}) > 0, \forall \mathbf{x} \in (0, \eta)$.

081
082 **Lemma 1** [2] Let Ω be a compact set and let $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ be a proper and lower
083 semicontinuous function. Assume that f is constant on Ω and satisfies the KL property at each point
084 of Ω . Then there exists $\epsilon > 0, \eta > 0$ and $\varphi \in \Phi_\eta$, such that for all $\bar{\mathbf{u}}$ in Ω and all \mathbf{u} in the following
085 intersection

$$086 \{\mathbf{u} \in \mathbb{R}^n : \operatorname{dist}(\mathbf{u}, \Omega) < \epsilon\} \cap \{\mathbf{u} \in \mathbb{R}^n : f(\bar{\mathbf{u}}) < f(\mathbf{u}) < f(\bar{\mathbf{u}}) + \eta\}, \quad (10)$$

087
088 the following inequality holds

$$089 \varphi'(f(\mathbf{u}) - f(\bar{\mathbf{u}})) \operatorname{dist}(0, \partial f(\mathbf{u})) > 1, \quad (11)$$

090
091 All semi-algebraic functions and subanalytic functions satisfy the KL property [3, 2]. So KL prop-
092 erty is general enough. Typical examples include: real polynomial functions, logistic loss function
093 $\log(1 + e^{-t})$, $\|\mathbf{x}\|_p$ ($p \geq 0$), $\|\mathbf{x}\|_\infty$, indicator function of the positive semidefinite (PSD) cone, the
094 Stiefel manifolds and the set of constant rank matrices.

095 096 2 Monotone APG

097
098 We summarize the monotone APG in Algorithm 1 and monotone APG with line search in Algorithm
099 2.

100
101 **Theorem 1** Let f be a proper function with Lipschitz continuous gradients and g be proper and
102 lower semicontinuous. For nonconvex f and nonconvex nonsmooth g , assume that (3) holds. Then
103 $\{\mathbf{x}_k\}$ and $\{\mathbf{v}_k\}$ generated by Algorithm 1 are bounded. Let \mathbf{x}^* be any accumulation point of $\{\mathbf{x}_k\}$,
104 we have $0 \in \partial F(\mathbf{x}^*)$.

Algorithm 1 monotone APG with fixed stepsize

Initialize $\mathbf{z}_1 = \mathbf{x}_1 = \mathbf{x}_0$, $t_1 = 1$, $t_0 = 0$, $\alpha_y < \frac{1}{L}$, $\alpha_x < \frac{1}{L}$.
for $k = 1, 2, 3, \dots$ **do**

$$\mathbf{y}_k = \mathbf{x}_k + \frac{t_{k-1}}{t_k}(\mathbf{z}_k - \mathbf{x}_k) + \frac{t_{k-1} - 1}{t_k}(\mathbf{x}_k - \mathbf{x}_{k-1}), \quad (12)$$

$$\mathbf{z}_{k+1} = \text{prox}_{\alpha_y g}(\mathbf{y}_k - \alpha_y \nabla f(\mathbf{y}_k)), \quad (13)$$

$$\mathbf{v}_{k+1} = \text{prox}_{\alpha_x g}(\mathbf{x}_k - \alpha_x \nabla f(\mathbf{x}_k)), \quad (14)$$

$$t_{k+1} = \frac{\sqrt{4(t_k)^2 + 1} + 1}{2}, \quad (15)$$

$$\mathbf{x}_{k+1} = \begin{cases} \mathbf{z}_{k+1}, & \text{if } F(\mathbf{z}_{k+1}) \leq F(\mathbf{v}_{k+1}), \\ \mathbf{v}_{k+1}, & \text{otherwise.} \end{cases} \quad (16)$$

end for

Proof (14) in Algorithm 1 can be seen as

$$\mathbf{v}_{k+1} = \underset{\mathbf{x}}{\text{argmin}} \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\alpha_x} \|\mathbf{x} - \mathbf{x}_k\|^2 + g(\mathbf{x}). \quad (17)$$

So we have

$$\langle \nabla f(\mathbf{x}_k), \mathbf{v}_{k+1} - \mathbf{x}_k \rangle + \frac{1}{2\alpha_x} \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 + g(\mathbf{v}_{k+1}) \leq g(\mathbf{x}_k). \quad (18)$$

From the Lipschitz continuous of ∇f we have

$$F(\mathbf{v}_{k+1}) \leq g(\mathbf{v}_{k+1}) + f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{v}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \quad (19)$$

$$\leq g(\mathbf{x}_k) - \langle \nabla f(\mathbf{x}_k), \mathbf{v}_{k+1} - \mathbf{x}_k \rangle - \frac{1}{2\alpha_x} \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \quad (20)$$

$$+ f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{v}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \quad (21)$$

$$= F(\mathbf{x}_k) - \left(\frac{1}{2\alpha_x} - \frac{L}{2} \right) \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2. \quad (22)$$

If $F(\mathbf{z}_{k+1}) \leq F(\mathbf{v}_{k+1})$, then

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1}, F(\mathbf{x}_{k+1}) = F(\mathbf{z}_{k+1}) \leq F(\mathbf{v}_{k+1}). \quad (23)$$

If $F(\mathbf{z}_{k+1}) > F(\mathbf{v}_{k+1})$, then

$$\mathbf{x}_{k+1} = \mathbf{v}_{k+1}, F(\mathbf{x}_{k+1}) = F(\mathbf{v}_{k+1}). \quad (24)$$

From (22), (23) and (24) we have

$$F(\mathbf{x}_{k+1}) \leq F(\mathbf{v}_{k+1}) \leq F(\mathbf{x}_k). \quad (25)$$

So

$$F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_1), F(\mathbf{v}_{k+1}) \leq F(\mathbf{x}_1) \quad (26)$$

for all k . From the assumption we know that $\{\mathbf{x}_k\}$ and $\{\mathbf{v}_k\}$ are bounded. Thus $\{\mathbf{x}_k\}$ has accumulation points. As $F(\mathbf{x}_k)$ is nonincreasing, F has the same value at all the accumulation points. Let it be F^* . From (22) we have

$$\left(\frac{1}{2\alpha_x} - \frac{L}{2} \right) \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \leq F(\mathbf{x}_k) - F(\mathbf{v}_{k+1}) \leq F(\mathbf{x}_k) - F(\mathbf{x}_{k+1}). \quad (27)$$

Summing over $k = 1, 2, \dots, \infty$, we have

$$\left(\frac{1}{2\alpha_x} - \frac{L}{2} \right) \sum_{k=1}^{\infty} \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \leq F(\mathbf{x}_1) - F^* < \infty, \quad (28)$$

From $\alpha_x < \frac{1}{L}$ we have

$$\|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \rightarrow 0 \quad \text{as } k \rightarrow \infty \quad (29)$$

From the optimality condition of (17) we have

$$0 \in \nabla f(\mathbf{x}_k) + \frac{1}{\alpha_x}(\mathbf{v}_{k+1} - \mathbf{x}_k) + \partial g(\mathbf{v}_{k+1}) \quad (30)$$

$$= \nabla f(\mathbf{v}_{k+1}) + \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{v}_{k+1}) + \frac{1}{\alpha_x}(\mathbf{v}_{k+1} - \mathbf{x}_k) + \partial g(\mathbf{v}_{k+1}). \quad (31)$$

So we have

$$-\nabla f(\mathbf{x}_k) + \nabla f(\mathbf{v}_{k+1}) - \frac{1}{\alpha_x}(\mathbf{v}_{k+1} - \mathbf{x}_k) \in \partial F(\mathbf{v}_{k+1}), \quad (32)$$

and

$$\left\| \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{v}_{k+1}) + \frac{1}{\alpha_x}(\mathbf{v}_{k+1} - \mathbf{x}_k) \right\| \leq \left(\frac{1}{\alpha_x} + L \right) \|\mathbf{v}_{k+1} - \mathbf{x}_k\| \rightarrow 0, \quad (33)$$

as $k \rightarrow \infty$.

Let \mathbf{x}^* be any accumulation point of $\{\mathbf{x}_k\}$, say $\{\mathbf{x}_{k_j}\} \rightarrow \mathbf{x}^*$ as $j \rightarrow \infty$. From (29) we have $\{\mathbf{v}_{k_j+1}\} \rightarrow \mathbf{x}^*$ as $j \rightarrow \infty$. From (17) we have

$$\langle \nabla f(\mathbf{x}_{k_j}), \mathbf{v}_{k_j+1} - \mathbf{x}_{k_j} \rangle + \frac{1}{2\alpha_x} \|\mathbf{v}_{k_j+1} - \mathbf{x}_{k_j}\|^2 + g(\mathbf{v}_{k_j+1}) \quad (34)$$

$$\leq \langle \nabla f(\mathbf{x}_{k_j}), \mathbf{x}^* - \mathbf{x}_{k_j} \rangle + \frac{1}{2\alpha_x} \|\mathbf{x}^* - \mathbf{x}_{k_j}\|^2 + g(\mathbf{x}^*). \quad (35)$$

So

$$\limsup_{j \rightarrow \infty} g(\mathbf{v}_{k_j+1}) \leq g(\mathbf{x}^*). \quad (36)$$

From the definition of lower semicontinuous of g we have

$$\liminf_{j \rightarrow \infty} g(\mathbf{v}_{k_j+1}) \geq g(\mathbf{x}^*). \quad (37)$$

So we have

$$\lim_{j \rightarrow \infty} g(\mathbf{v}_{k_j+1}) = g(\mathbf{x}^*). \quad (38)$$

Because f is continuously differentiable, we have

$$\lim_{j \rightarrow \infty} F(\mathbf{v}_{k_j+1}) = F(\mathbf{x}^*). \quad (39)$$

From $\{\mathbf{v}_{k_j+1}\} \rightarrow \mathbf{x}^*$, (39), (32), (33) and Proposition 1.2 we have

$$0 \in \partial F(\mathbf{x}^*). \quad (40)$$

■

Corollary 1 *Let f be a proper function with Lipschitz continuous gradients and g be proper and lower semicontinuous. For nonconvex f and nonconvex nonsmooth g , assume that (3) holds, then $\{\mathbf{x}_k\}$ and $\{\mathbf{v}_k\}$ generated by Algorithm 2 are bounded. Let \mathbf{x}^* be any accumulation point of $\{\mathbf{x}_k\}$, we have $0 \in \partial F(\mathbf{x}^*)$.*

Proof From (22) and similar deduction we know that such α_y and α_x satisfying

$$\mathbf{v}_{k+1} = \text{prox}_{\alpha_x g}(\mathbf{x}_k - \alpha_x \nabla f(\mathbf{x}_k)), \quad (52)$$

$$F(\mathbf{v}_{k+1}) \leq F(\mathbf{x}_k) - \delta \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2, \quad (53)$$

$$\mathbf{z}_{k+1} = \text{prox}_{\alpha_y g}(\mathbf{y}_k - \alpha_y \nabla f(\mathbf{y}_k)), \quad (54)$$

$$F(\mathbf{z}_{k+1}) \leq F(\mathbf{y}_k) - \delta \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2, \quad (55)$$

exist, e.g., when they are reduced until $\alpha_x < \frac{1}{L}$ and $\alpha_y < \frac{1}{L}$. So the line search can be terminated in finite iterations. Similar to Theorem 1 we can have the conclusion. ■

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Algorithm 2 monotone APG with line search

Initialize $\mathbf{z}_1 = \mathbf{x}_1 = \mathbf{x}_0$, $t_1 = 1$, $t_0 = 0$, $\delta > 0$, $\rho < 1$.
for $k = 1, 2, 3, \dots$ **do**

$$\mathbf{y}_k = \mathbf{x}_k + \frac{t_{k-1}}{t_k}(\mathbf{z}_k - \mathbf{x}_k) + \frac{t_{k-1} - 1}{t_k}(\mathbf{x}_k - \mathbf{x}_{k-1}), \quad (41)$$

$$\mathbf{s}_k = \mathbf{z}_k - \mathbf{y}_{k-1}, \mathbf{r}_k = \nabla f(\mathbf{z}_k) - \nabla f(\mathbf{y}_{k-1}), \quad (42)$$

$$\alpha_y = \frac{(\mathbf{s}_k)^T \mathbf{s}_k}{(\mathbf{s}_k)^T \mathbf{r}_k} \quad \text{or} \quad \alpha_y = \frac{(\mathbf{s}_k)^T \mathbf{r}_k}{(\mathbf{r}_k)^T \mathbf{r}_k}, \quad (43)$$

$$\mathbf{s}_k = \mathbf{v}_k - \mathbf{x}_{k-1}, \mathbf{r}_k = \nabla f(\mathbf{v}_k) - \nabla f(\mathbf{x}_{k-1}), \quad (44)$$

$$\alpha_x = \frac{(\mathbf{s}_k)^T \mathbf{s}_k}{(\mathbf{s}_k)^T \mathbf{r}_k} \quad \text{or} \quad \alpha_x = \frac{(\mathbf{s}_k)^T \mathbf{r}_k}{(\mathbf{r}_k)^T \mathbf{r}_k}. \quad (45)$$

Repeat

$$\mathbf{z}_{k+1} = \text{prox}_{\alpha_y g}(\mathbf{y}_k - \alpha_y \nabla f(\mathbf{y}_k)), \quad (46)$$

$$\alpha_y = \alpha_y \times \rho, \quad (47)$$

until $F(\mathbf{z}_{k+1}) \leq F(\mathbf{y}_k) - \delta \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2$.

Repeat

$$\mathbf{v}_{k+1} = \text{prox}_{\alpha_x g}(\mathbf{x}_k - \alpha_x \nabla f(\mathbf{x}_k)), \quad (48)$$

$$\alpha_x = \alpha_x \times \rho, \quad (49)$$

until $F(\mathbf{v}_{k+1}) \leq F(\mathbf{x}_k) - \delta \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2$.

$$t_{k+1} = \frac{\sqrt{4(t_k)^2 + 1} + 1}{2}, \quad (50)$$

$$\mathbf{x}_{k+1} = \begin{cases} \mathbf{z}_{k+1}, & \text{if } F(\mathbf{z}_{k+1}) \leq F(\mathbf{v}_{k+1}), \\ \mathbf{v}_{k+1}, & \text{otherwise.} \end{cases} \quad (51)$$

end for

Theorem 2 Assume that f and g are convex and ∇f is Lipschitz continuous. Then $\{\mathbf{x}_k\}$ generated by algorithm 1 satisfies

$$F(\mathbf{x}_{N+1}) - F(\mathbf{x}^*) \leq \frac{2}{\alpha_y(N+1)^2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \quad (56)$$

where \mathbf{x}^* is a global minimizer of $F(\mathbf{x})$.

Proof (13) in Algorithm 1 can be seen as

$$\mathbf{z}_{k+1} = \underset{\mathbf{x}}{\text{argmin}} \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{1}{2\alpha_y} \|\mathbf{x} - \mathbf{y}_k\|^2 + g(\mathbf{x}). \quad (57)$$

From the optimality condition, we have

$$0 \in \nabla f(\mathbf{y}_k) + \frac{1}{\alpha_y}(\mathbf{z}_{k+1} - \mathbf{y}_k) + \partial g(\mathbf{z}_{k+1}). \quad (58)$$

From the convexity of g we have

$$g(\mathbf{x}) - g(\mathbf{z}_{k+1}) \geq \left\langle -\nabla f(\mathbf{y}_k) - \frac{1}{\alpha_y}(\mathbf{z}_{k+1} - \mathbf{y}_k), \mathbf{x} - \mathbf{z}_{k+1} \right\rangle, \forall \mathbf{x}. \quad (59)$$

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

From the Lipschitz continuous of ∇f and convexity of f we have

$$F(\mathbf{z}_{k+1}) \leq g(\mathbf{z}_{k+1}) + f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{z}_{k+1} - \mathbf{y}_k \rangle + \frac{L}{2} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2 \quad (60)$$

$$= g(\mathbf{z}_{k+1}) + f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \langle \nabla f(\mathbf{y}_k), \mathbf{z}_{k+1} - \mathbf{x} \rangle \quad (61)$$

$$+ \frac{L}{2} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2 \quad (62)$$

$$\leq g(\mathbf{z}_{k+1}) + f(\mathbf{x}) + \langle \nabla f(\mathbf{y}_k), \mathbf{z}_{k+1} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2 \quad (63)$$

$$\leq g(\mathbf{x}) + \left\langle \nabla f(\mathbf{y}_k) + \frac{1}{\alpha_y} (\mathbf{z}_{k+1} - \mathbf{y}_k), \mathbf{x} - \mathbf{z}_{k+1} \right\rangle \quad (64)$$

$$+ f(\mathbf{x}) + \langle \nabla f(\mathbf{y}_k), \mathbf{z}_{k+1} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2 \quad (65)$$

$$= F(\mathbf{x}) + \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, \mathbf{x} - \mathbf{z}_{k+1} \rangle + \frac{L}{2} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2 \quad (66)$$

$$= F(\mathbf{x}) + \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, \mathbf{x} - \mathbf{y}_k + \mathbf{y}_k - \mathbf{z}_{k+1} \rangle + \frac{L}{2} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2 \quad (67)$$

$$= F(\mathbf{x}) + \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, \mathbf{x} - \mathbf{y}_k \rangle - \left(\frac{1}{\alpha_y} - \frac{L}{2} \right) \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2 \quad (68)$$

$$\leq F(\mathbf{x}) + \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, \mathbf{x} - \mathbf{y}_k \rangle - \frac{1}{2\alpha_y} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2. \quad (69)$$

Let $\mathbf{x} = \mathbf{x}_k$ and \mathbf{x}^* , we have

$$F(\mathbf{z}_{k+1}) - F(\mathbf{x}_k) \leq \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, \mathbf{x}_k - \mathbf{y}_k \rangle - \frac{1}{2\alpha_y} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2, \quad (70)$$

$$F(\mathbf{z}_{k+1}) - F(\mathbf{x}^*) \leq \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, \mathbf{x}^* - \mathbf{y}_k \rangle - \frac{1}{2\alpha_y} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2. \quad (71)$$

Multiplying (70) by $t_k - 1$ and adding (71) we have

$$t_k F(\mathbf{z}_{k+1}) - (t_k - 1)F(\mathbf{x}_k) - F(\mathbf{x}^*) \quad (72)$$

$$\leq \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, (t_k - 1)(\mathbf{x}_k - \mathbf{y}_k) + \mathbf{x}^* - \mathbf{y}_k \rangle - \frac{t_k}{2\alpha_y} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2. \quad (73)$$

So we have

$$t_k (F(\mathbf{z}_{k+1}) - F(\mathbf{x}^*)) - (t_k - 1) (F(\mathbf{x}_k) - F(\mathbf{x}^*)) \quad (74)$$

$$\leq \frac{1}{\alpha_y} \langle \mathbf{z}_{k+1} - \mathbf{y}_k, (t_k - 1)(\mathbf{x}_k - \mathbf{y}_k) + \mathbf{x}^* - \mathbf{y}_k \rangle - \frac{t_k}{2\alpha_y} \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2. \quad (75)$$

Multiplying both sides by t_k and using $(t_k)^2 - t_k = (t_{k-1})^2$ from (15) we have

$$(t_k)^2 (F(\mathbf{z}_{k+1}) - F(\mathbf{x}^*)) - (t_{k-1})^2 (F(\mathbf{x}_k) - F(\mathbf{x}^*)) \quad (76)$$

$$\leq \frac{1}{\alpha_y} \langle t_k (\mathbf{z}_{k+1} - \mathbf{y}_k), (t_k - 1)(\mathbf{x}_k - \mathbf{y}_k) + \mathbf{x}^* - \mathbf{y}_k \rangle - \frac{1}{2\alpha_y} \|t_k (\mathbf{z}_{k+1} - \mathbf{y}_k)\|^2 \quad (77)$$

$$= \frac{1}{\alpha_y} \langle t_k (\mathbf{z}_{k+1} - \mathbf{y}_k), (t_k - 1)\mathbf{x}_k - t_k \mathbf{y}_k + \mathbf{x}^* \rangle - \frac{1}{2\alpha_y} \|t_k (\mathbf{z}_{k+1} - \mathbf{y}_k)\|^2 \quad (78)$$

$$= \frac{1}{2\alpha_y} (\|(t_k - 1)\mathbf{x}_k - t_k \mathbf{y}_k + \mathbf{x}^*\|^2 - \|(t_k - 1)\mathbf{x}_k - t_k \mathbf{z}_{k+1} + \mathbf{x}^*\|^2). \quad (79)$$

Define

$$U_{k+1} = t_k \mathbf{z}_{k+1} - (t_k - 1)\mathbf{x}_k - \mathbf{x}^*. \quad (80)$$

Let

$$U_k = t_{k-1} \mathbf{z}_k - (t_{k-1} - 1)\mathbf{x}_{k-1} - \mathbf{x}^* = t_k \mathbf{y}_k - (t_k - 1)\mathbf{x}_k - \mathbf{x}^*. \quad (81)$$

324 We have
325

$$326 \quad \mathbf{y}_k = \frac{t_{k-1}\mathbf{z}_k - (t_{k-1} - 1)\mathbf{x}_{k-1} + (t_k - 1)\mathbf{x}_k}{t_k} \quad (82)$$

$$327 \quad = \mathbf{x}_k + \frac{t_{k-1}}{t_k}(\mathbf{z}_k - \mathbf{x}_k) + \frac{t_{k-1} - 1}{t_k}(\mathbf{x}_k - \mathbf{x}_{k-1}), \quad (83)$$

330 which is the same with (12) in Algorithm 1. So we have
331

$$332 \quad (t_k)^2 (F(\mathbf{z}_{k+1}) - F(\mathbf{x}^*)) - (t_{k-1})^2 (F(\mathbf{x}_k) - F(\mathbf{x}^*)) \quad (84)$$

$$333 \quad \leq \frac{1}{2\alpha_y} (\|U_k\|^2 - \|U_{k+1}\|^2). \quad (85)$$

336 If $F(\mathbf{z}_{k+1}) \leq F(\mathbf{v}_{k+1})$, then $\mathbf{x}_{k+1} = \mathbf{z}_{k+1}$. So
337

$$338 \quad (t_k)^2 (F(\mathbf{x}_{k+1}) - F(\mathbf{x}^*)) - (t_{k-1})^2 (F(\mathbf{x}_k) - F(\mathbf{x}^*)) \quad (86)$$

$$339 \quad = (t_k)^2 (F(\mathbf{z}_{k+1}) - F(\mathbf{x}^*)) - (t_{k-1})^2 (F(\mathbf{x}_k) - F(\mathbf{x}^*)) \quad (87)$$

$$340 \quad \leq \frac{1}{2\alpha_y} (\|U_k\|^2 - \|U_{k+1}\|^2). \quad (88)$$

343 If $F(\mathbf{z}_{k+1}) > F(\mathbf{v}_{k+1})$, then $\mathbf{x}_{k+1} = \mathbf{v}_{k+1}$. So
344

$$345 \quad (t_k)^2 (F(\mathbf{x}_{k+1}) - F(\mathbf{x}^*)) - (t_{k-1})^2 (F(\mathbf{x}_k) - F(\mathbf{x}^*)) \quad (89)$$

$$346 \quad \leq (t_k)^2 (F(\mathbf{z}_{k+1}) - F(\mathbf{x}^*)) - (t_{k-1})^2 (F(\mathbf{x}_k) - F(\mathbf{x}^*)) \quad (90)$$

$$347 \quad \leq \frac{1}{2\alpha_y} (\|U_k\|^2 - \|U_{k+1}\|^2). \quad (91)$$

349 Summing over $k = 1, \dots, N$, we have
350

$$351 \quad (t_N)^2 (F(\mathbf{x}_{N+1}) - F(\mathbf{x}^*)) \quad (92)$$

$$352 \quad = (t_N)^2 (F(\mathbf{x}_{N+1}) - F(\mathbf{x}^*)) - (t^0)^2 (F(\mathbf{x}_1) - F(\mathbf{x}^*)) \quad (93)$$

$$353 \quad \leq \frac{1}{2\alpha_y} (\|U_1\|^2 - \|U_{N+1}\|^2) \quad (94)$$

$$354 \quad \leq \frac{1}{2\alpha_y} \|U_1\|^2 \quad (95)$$

$$355 \quad = \frac{1}{2\alpha_y} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \quad (96)$$

361 From (15) we can easily have that $t_k \geq \frac{k+1}{2}$. So we have
362

$$363 \quad F(\mathbf{x}_{N+1}) - F(\mathbf{x}^*) \leq \frac{2}{\alpha_y(N+1)^2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \quad (97)$$

366 ■

367 **Theorem 3** Let f be a proper function with Lipschitz continuous gradients and g be proper and
368 lower semicontinuous. For nonconvex f and nonconvex nonsmooth g , assume that (3) holds. If we
369 further assume that f and g satisfy the KL property, and the desingularising function has the form
370 of $\varphi(t) = \frac{C}{\theta} t^\theta$ for some $C > 0$, $\theta \in (0, 1]$, then
371

372 1. If $\theta = 1$, then there exists k_1 such that $F(\mathbf{x}_k) = F^*$ for all $k > k_1$ and the algorithm
373 terminates in finite steps.

374 2. If $\theta \in [\frac{1}{2}, 1)$, then there exists k_2 such that for all $k > k_2$,

$$375 \quad F(\mathbf{x}_k) - F^* \leq \left(\frac{d_1 C^2}{1 + d_1 C^2} \right)^{k-k_2} r_{k_2}. \quad (98)$$

378 3. If $\theta \in (0, \frac{1}{2})$, then there exists k_3 such that for all $k > k_3$,
 379

$$380 \quad F(\mathbf{x}_k) - F^* \leq \left(\frac{C}{(k - k_3)d_2(1 - 2\theta)} \right)^{\frac{1}{1-2\theta}}, \quad (99)$$

382 where F^* is the same function value at all the accumulation points of $\{\mathbf{x}_k\}$, $r_k = F(\mathbf{v}_k) -$
 383 F^* , $d_1 = \left(\frac{1}{\alpha_x} + L\right)^2 / \left(\frac{1}{2\alpha_x} - \frac{L}{2}\right)$, $d_2 = \min \left\{ \frac{1}{2d_1C}, \frac{C}{1-2\theta} \left(2^{\frac{2\theta-1}{2\theta-2}} - 1\right) r_0^{2\theta-1} \right\}$
 384

385 **Proof** From (22) and (25) we have
 386

$$387 \quad F(\mathbf{v}_{k+1}) \leq F(\mathbf{x}_k) - \left(\frac{1}{2\alpha_x} - \frac{L}{2}\right) \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \quad (100)$$

$$388 \quad \leq F(\mathbf{v}_k) - \left(\frac{1}{2\alpha_x} - \frac{L}{2}\right) \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2. \quad (101)$$

389 From (33) we have
 390

$$391 \quad \text{dist}(0, \partial F(\mathbf{v}_{k+1})) \leq \left(\frac{1}{\alpha_x} + L\right) \|\mathbf{v}_{k+1} - \mathbf{x}_k\|. \quad (102)$$

392 From (29) we know that $\{\mathbf{x}_k\}$ and $\{\mathbf{v}_k\}$ have the same accumulation points. Let Ω be the set
 393 that contains all the accumulation points of $\{\mathbf{x}_k\}$ (also $\{\mathbf{v}_k\}$). Because $F(\mathbf{v}_k)$ is nonincreasing, F
 394 has the same value at all the accumulation points in Ω . Let it be F^* . So we have
 395

$$396 \quad F(\mathbf{v}_k) \geq F^*, F(\mathbf{v}_k) \rightarrow F^*. \quad (103)$$

397 If there exists \bar{k} such that $F(\mathbf{v}^{\bar{k}}) = F^*$, then $F(\mathbf{v}^{\bar{k}}) = F(\mathbf{v}^{\bar{k}+1}) = \dots = F^*$. So $\|\mathbf{v}^{\bar{k}+1} -$
 398 $\mathbf{x}^{\bar{k}}\| = \|\mathbf{v}^{\bar{k}+2} - \mathbf{x}^{\bar{k}+1}\| = \dots = 0$. The conclusion holds. If $F(\mathbf{v}_k) > F^*$ for all k , then from
 399 $F(\mathbf{v}_k) \rightarrow F^*$ we know that there exists \hat{k}_1 such that $F(\mathbf{v}_k) < F^* + \eta$ whenever $k > \hat{k}_1$. On the
 400 other hand, because $\text{dist}(\mathbf{v}_k, \Omega) \rightarrow 0$, there exists \hat{k}_2 such that $\text{dist}(\mathbf{v}_k, \Omega) < \varepsilon$ whenever $k > \hat{k}_2$.
 401 Let $k > k_0 = \max\{\hat{k}_1, \hat{k}_2\}$, we have
 402

$$403 \quad \mathbf{v}_k \in \{\mathbf{v}, \text{dist}(\mathbf{v}, \Omega) \leq \varepsilon\} \cap [F^* < F(\mathbf{v}) < F^* + \eta]. \quad (104)$$

404 From the uniform KL property in Lemma 1, there exists a concave function φ such that
 405

$$406 \quad \varphi'(F(\mathbf{v}_k) - F^*) \text{dist}(0, \partial F(\mathbf{v}_k)) \geq 1. \quad (105)$$

407 Define $r_k = F(\mathbf{v}_k) - F^*$. We suppose that $r_k > 0$ for all k . Otherwise $F(\mathbf{v}_k) = F(\mathbf{v}_{k+1}) = \dots =$
 408 F^* and the algorithm terminates in finite steps. By supposing this (105) holds.
 409

410 From (102), (105) and (101) we have
 411

$$412 \quad 1 \leq [\varphi'(F(\mathbf{v}_k) - F^*) \text{dist}(0, \partial F(\mathbf{v}_k))]^2 \quad (106)$$

$$413 \quad \leq [\varphi'(r_k)]^2 \left(\frac{1}{\alpha_x} + L\right)^2 \|\mathbf{v}_k - \mathbf{x}_{k-1}\|^2 \quad (107)$$

$$414 \quad \leq [\varphi'(r_k)]^2 \left(\frac{1}{\alpha_x} + L\right)^2 \frac{F(\mathbf{v}_{k-1}) - F(\mathbf{v}_k)}{\left(\frac{1}{2\alpha_x} - \frac{L}{2}\right)} \quad (108)$$

$$415 \quad = d_1 [\varphi'(r_k)]^2 (r_{k-1} - r_k), \quad (109)$$

416 for all $k > k_0$, where $d_1 = \left(\frac{1}{\alpha_x} + L\right)^2 / \left(\frac{1}{2\alpha_x} - \frac{L}{2}\right)$. Because φ has the form of $\varphi(t) = \frac{C}{\theta} t^\theta$, we
 417 have $\varphi'(t) = Ct^{\theta-1}$. So (109) becomes
 418

$$419 \quad 1 \leq d_1 C^2 r_k^{2\theta-2} (r_{k-1} - r_k). \quad (110)$$

420 1. Case $\theta = 1$.
 421

422 In this case, (110) becomes
 423

$$424 \quad 1 \leq d_1 C^2 (r_k - r_{k+1}). \quad (111)$$

432 Because $r_k \rightarrow 0$ and $d_1 > 0, C > 0$, this is a contradiction. So there exists k_1 such that $r_k = 0$ for
 433 all $k > k_1$. The algorithm terminates in finite steps.

434 2. Case $\theta \in [\frac{1}{2}, 1)$.

435 In this case, $0 < 2 - 2\theta \leq 1$. As $r_k \rightarrow 0$, there exists \hat{k}_3 such that $r_k^{2-2\theta} \geq r_k$ for all $k > \hat{k}_3$.
 436 (110) becomes

$$437 r_k \leq d_1 C^2 (r_{k-1} - r_k). \quad (112)$$

438 So we have

$$439 r_k \leq \frac{d_1 C^2}{1 + d_1 C^2} r_{k-1}, \quad (113)$$

440 for all $k_2 > \max\{k_0, \hat{k}_3\}$ and

$$441 r_k \leq \left(\frac{d_1 C^2}{1 + d_1 C^2} \right)^{k-k_2} r_{k_2}. \quad (114)$$

442 So we have

$$443 F(\mathbf{x}_k) - F^* \leq F(\mathbf{v}_k) - F^* = r_k \leq \left(\frac{d_1 C^2}{1 + d_1 C^2} \right)^{k-k_2} r_{k_2}. \quad (115)$$

444 3. Case $\theta \in (0, \frac{1}{2})$.

445 In this case, $2\theta - 2 \in (-2, -1), 2\theta - 1 \in (-1, 0)$. As $r_{k-1} > r_k$, we have $r_{k-1}^{2\theta-2} < r_k^{2\theta-2}$ and
 446 $r_0^{2\theta-1} < \dots < r_{k-1}^{2\theta-1} < r_k^{2\theta-1}$

447 Define $\phi(t) = \frac{C}{1-2\theta} t^{2\theta-1}$, then $\phi'(t) = -Ct^{2\theta-2}$.

448 If $r_k^{2\theta-2} \leq 2r_{k-1}^{2\theta-2}$, then

$$449 \phi(r_k) - \phi(r_{k-1}) = \int_{r_{k-1}}^{r_k} \phi'(t) dt = C \int_{r_k}^{r_{k-1}} t^{2\theta-2} dt \quad (116)$$

$$450 \geq C(r_{k-1} - r_k) r_{k-1}^{2\theta-2} \geq \frac{C}{2} (r_{k-1} - r_k) r_k^{2\theta-2} \quad (117)$$

$$451 \geq \frac{1}{2d_1 C}. \quad (118)$$

452 for all $k > k_0$.

453 If $r_k^{2\theta-2} \geq 2r_{k-1}^{2\theta-2}$, then $r_k^{2\theta-1} \geq 2^{\frac{2\theta-1}{2\theta-2}} r_{k-1}^{2\theta-1}$.

$$454 \phi(r_k) - \phi(r_{k-1}) = \frac{C}{1-2\theta} (r_k^{2\theta-1} - r_{k-1}^{2\theta-1}) \quad (119)$$

$$455 \geq \frac{C}{1-2\theta} (2^{\frac{2\theta-1}{2\theta-2}} - 1) r_{k-1}^{2\theta-1} \quad (120)$$

$$456 = q r_{k-1}^{2\theta-1} \geq q r_0^{2\theta-1}. \quad (121)$$

457 where $q = \frac{C}{1-2\theta} (2^{\frac{2\theta-1}{2\theta-2}} - 1)$. Let $d_2 = \min\{\frac{1}{2d_1 C}, q r_0^{2\theta-1}\}$, we have

$$458 \phi(r_k) - \phi(r_{k-1}) \geq d_2, \quad (122)$$

459 for all $k > k_0$ and

$$460 \phi(r_k) \geq \phi(r_k) - \phi(r_{k_0}) \geq \sum_{i=k_0+1}^k \phi(r_i) - \phi(r_{i-1}) \geq (k - k_0) d_2. \quad (123)$$

461 So we have

$$462 r_k^{2\theta-1} \geq \frac{(k - k_0) d_2 (1 - 2\theta)}{C}, \quad (124)$$

486 and
487

$$488 \quad r_k \leq \left(\frac{C}{(k - k_0)d_2(1 - 2\theta)} \right)^{\frac{1}{1-2\theta}}. \quad (125)$$

489 Let $k_3 = k_0$ we have
491

$$492 \quad F(\mathbf{x}_k) - F^* \leq F(\mathbf{v}_k) - F^* = r_k \leq \left(\frac{C}{(k - k_3)d_2(1 - 2\theta)} \right)^{\frac{1}{1-2\theta}}, \quad (126)$$

493 which completes the proof. ■
494

495 Difference with the conditions in [5]:

496 [5] considered general descent method with the conditions:

$$497 \quad F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_k) - \alpha \|\mathbf{x}_{k+1} - \mathbf{x}_k\|, \quad (127)$$

500 and
501

$$502 \quad \|\partial F(\mathbf{x}_{k+1})\| \leq \beta \|\mathbf{x}_{k+1} - \mathbf{x}_k\| \quad (128)$$

503 Proximal gradient method is a typical example satisfying these conditions. However, to make the
504 proximal gradient method both accelerate and converge, we introduce the intermediate variables \mathbf{y}_k ,
505 \mathbf{v}_k and \mathbf{z}_k . This makes our algorithm more complex and the conditions satisfied by our algorithm
506 becomes

$$507 \quad F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_k) - \alpha \|\mathbf{v}_{k+1} - \mathbf{x}_k\|, F(\mathbf{v}_{k+1}) \leq F(\mathbf{v}_k) - \alpha \|\mathbf{v}_{k+1} - \mathbf{x}_k\| \quad (129)$$

508 and
509

$$510 \quad \|\partial F(\mathbf{v}_{k+1})\| \leq \beta \|\mathbf{v}_{k+1} - \mathbf{x}_k\| \quad (130)$$

511 The intermediate variable \mathbf{v} makes the main difference. As a result, under the conditions of (127)
512 and (128), a useful conclusion of finite length of $\{\mathbf{x}\}$: $\sum_{i=k}^{\infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\| < \infty$ can be achieved and
513 $\{\mathbf{x}_k\}$ is a converged sequence. Accordingly, the convergence rate for $\|\mathbf{x}_k - \mathbf{x}^*\|$ can be obtained.
514 By contrast, our algorithm can only get $\sum_{i=1}^{\infty} \|\mathbf{v}_{k+1} - \mathbf{x}_k\| < \infty$. Neither the convergence rate for
515 $\|\mathbf{x}_k - \mathbf{x}^*\|$ nor $\{\mathbf{x}_k\}$ is a converged sequence can be obtained.
516

517 3 Nonmonotone APG

518

519 We summarize the nonmonotone APG in Algorithm 3 and nonmonotone APG with line search in
520 Algorithm 4.
521

522 **Lemma 2** *In Algorithms 3 and 4, we have*

523

$$524 \quad F(\mathbf{x}_k) \leq c_k \leq A_k, A_k = \frac{\sum_{i=1}^k F(\mathbf{x}_i)}{k}, \quad (153)$$

525

526 and there exists α_x such that

527

$$528 \quad \mathbf{v}_{k+1} = \text{prox}_{\alpha_x g}(\mathbf{x}_k - \alpha_x \nabla f(\mathbf{x}_k)) \quad (154)$$

529

530 satisfies

531

$$532 \quad F(\mathbf{v}_{k+1}) \leq c_k - \delta \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2, \quad (155)$$

533

534 where δ is any small positive constant.

535

536 **Proof** We prove by induction. For $k = 1$, $c_1 = F(\mathbf{x}_1)$. From (17)-(22) we know that $\alpha_x < \frac{1}{L}$
537 satisfies

538

$$539 \quad F(\mathbf{v}_2) \leq c_1 - \delta \|\mathbf{v}_2 - \mathbf{x}_1\|, \quad (156)$$

539

where

$$\mathbf{v}_2 = \text{prox}_{\alpha_x g}(\mathbf{x}_1 - \alpha_x \nabla f(\mathbf{x}_1)). \quad (157)$$

Algorithm 3 nonmonotone APG with fixed stepsize

Initialize $\mathbf{z}_1 = \mathbf{x}_1 = \mathbf{x}_0$, $t_1 = 1$, $t_0 = 0$, $\eta \in [0, 1)$, $\delta > 0$, $c_1 = F(\mathbf{x}_1)$, $q_1 = 1$, $\alpha_x < \frac{1}{L}$,
 $\alpha_y < \frac{1}{L}$.

for $k = 1, 2, 3, \dots$ **do**

$$\mathbf{y}_k = \mathbf{x}_k + \frac{t_{k-1}}{t_k}(\mathbf{z}_k - \mathbf{x}_k) + \frac{t_{k-1} - 1}{t_k}(\mathbf{x}_k - \mathbf{x}_{k-1}), \quad (131)$$

$$\mathbf{z}_{k+1} = \text{prox}_{\alpha_y g}(\mathbf{y}_k - \alpha_y \nabla f(\mathbf{y}_k)), \quad (132)$$

if $F(\mathbf{z}_{k+1}) \leq c_k - \delta \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2$ **then**

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1}. \quad (133)$$

else

$$\mathbf{v}_{k+1} = \text{prox}_{\alpha_x g}(\mathbf{x}_k - \alpha_x \nabla f(\mathbf{x}_k)), \quad (134)$$

$$\mathbf{x}_{k+1} = \begin{cases} \mathbf{z}_{k+1}, & \text{if } F(\mathbf{z}_{k+1}) \leq F(\mathbf{v}_{k+1}), \\ \mathbf{v}_{k+1}, & \text{otherwise.} \end{cases} \quad (135)$$

end if

$$t_{k+1} = \frac{\sqrt{4(t_k)^2 + 1} + 1}{2}, \quad (136)$$

$$q_{k+1} = \eta q_k + 1, \quad (137)$$

$$c_{k+1} = \frac{\eta q_k c_k + F(\mathbf{x}_{k+1})}{q_{k+1}}. \quad (138)$$

end for

If for all $k = 1, \dots, j$, the conclusions hold, then we consider $k = j + 1$. Define

$$D_{j+1}(t) = \frac{tc_j + F(\mathbf{x}_{j+1})}{t + 1}, \quad (158)$$

then

$$\frac{d}{dt} D_{j+1}(t) = \frac{c_j - F(\mathbf{x}_{j+1})}{(t + 1)^2}. \quad (159)$$

If (133) in Algorithm 3 (or (144) in Algorithm 4) is executed, then

$$F(\mathbf{x}_{j+1}) = F(\mathbf{z}_{j+1}) \leq c_j. \quad (160)$$

If (135) in Algorithm 3 (or (149) in Algorithm 4) is executed, by the induction step, we have that $F(\mathbf{v}_{j+1}) \leq c_j - \delta \|\mathbf{v}_{j+1} - \mathbf{x}_j\|$. So

$$F(\mathbf{x}_{j+1}) \leq F(\mathbf{v}_{j+1}) \leq c_j. \quad (161)$$

So we have

$$\frac{d}{dt} D_{j+1}(t) \geq 0, \quad (162)$$

which means that $D_{j+1}(t)$ is nondecreasing. So

$$F(\mathbf{x}_{j+1}) = D_{j+1}(0) \leq D_{j+1}(\eta q_j) = c_{j+1}. \quad (163)$$

From the definition of q_k we have

$$q_{k+1} = 1 + \sum_{i=1}^k \eta^i < k + 1, \quad (164)$$

Algorithm 4 nonmonotone APG with line search

Initialize $\mathbf{z}_1 = \mathbf{x}_1 = \mathbf{x}_0$, $t_1 = 1$, $t_0 = 0$, $\eta \in [0, 1)$, $\delta > 0$, $\rho < 1$, $c_1 = F(\mathbf{x}_1)$, $q_1 = 1$.
for $k = 1, 2, 3, \dots$ **do**

$$\mathbf{y}_k = \mathbf{x}_k + \frac{t_{k-1}}{t_k}(\mathbf{z}_k - \mathbf{x}_k) + \frac{t_{k-1} - 1}{t_k}(\mathbf{x}_k - \mathbf{x}_{k-1}), \quad (139)$$

$$\mathbf{s}_k = \mathbf{y}_k - \mathbf{y}_{k-1}, \mathbf{r}_k = \nabla f(\mathbf{y}_k) - \nabla f(\mathbf{y}_{k-1}), \quad (140)$$

$$\alpha_y = \frac{(\mathbf{s}_k)^T \mathbf{s}_k}{(\mathbf{s}_k)^T \mathbf{r}_k} \quad \text{or} \quad \alpha_y = \frac{(\mathbf{s}_k)^T \mathbf{r}_k}{(\mathbf{r}_k)^T \mathbf{r}_k}, \quad (141)$$

Repeat

$$\mathbf{z}_{k+1} = \text{prox}_{\alpha_y g}(\mathbf{y}_k - \alpha_y \nabla f(\mathbf{y}_k)), \quad (142)$$

$$\alpha_y = \alpha_y \times \rho, \quad (143)$$

until $F(\mathbf{z}_{k+1}) \leq F(\mathbf{y}_k) - \delta \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2$ or $F(\mathbf{z}_{k+1}) \leq c_k - \delta \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2$.
if $F(\mathbf{z}_{k+1}) \leq c_k - \delta \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2$ **then**

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1}. \quad (144)$$

else

$$\mathbf{s}_k = \mathbf{x}_k - \mathbf{y}_{k-1}, \mathbf{r}_k = \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{y}_{k-1}), \quad (145)$$

$$\alpha_x = \frac{(\mathbf{s}_k)^T \mathbf{s}_k}{(\mathbf{s}_k)^T \mathbf{r}_k} \quad \text{or} \quad \alpha_x = \frac{(\mathbf{s}_k)^T \mathbf{r}_k}{(\mathbf{r}_k)^T \mathbf{r}_k}, \quad (146)$$

Repeat

$$\mathbf{v}_{k+1} = \text{prox}_{\alpha_x g}(\mathbf{x}_k - \alpha_x \nabla f(\mathbf{x}_k)), \quad (147)$$

$$\alpha_x = \alpha_x \times \rho, \quad (148)$$

until $F(\mathbf{v}_{k+1}) \leq c_k - \delta \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2$.

$$\mathbf{x}_{k+1} = \begin{cases} \mathbf{z}_{k+1}, & \text{if } F(\mathbf{z}_{k+1}) \leq F(\mathbf{v}_{k+1}), \\ \mathbf{v}_{k+1}, & \text{otherwise.} \end{cases} \quad (149)$$

end if

$$t_{k+1} = \frac{\sqrt{4(t_k)^2 + 1} + 1}{2}, \quad (150)$$

$$q_{k+1} = \eta q_k + 1, \quad (151)$$

$$c_{k+1} = \frac{\eta q_k c_k + F(\mathbf{x}_{k+1})}{q_{k+1}}. \quad (152)$$

end for

due to $\eta \in [0, 1)$. So we have

$$c_{j+1} = D_{j+1}(\eta q_j) = D_{j+1}(q_{j+1} - 1) \quad (165)$$

$$\leq D_{j+1}(j) = \frac{j c_j + F(\mathbf{x}_{j+1})}{j+1} \leq \frac{j A_j + F(\mathbf{x}_{j+1})}{j+1} = A_{j+1}. \quad (166)$$

From (17)-(22) and using $F(\mathbf{x}_{j+1}) \leq c_{j+1}$ we have

$$F(\mathbf{v}_{j+2}) \leq F(\mathbf{x}_{j+1}) - \left(\frac{1}{2\alpha_x} - \frac{L}{2} \right) \|\mathbf{v}_{j+2} - \mathbf{x}_{j+1}\|^2 \quad (167)$$

$$\leq c_{j+1} - \left(\frac{1}{2\alpha_x} - \frac{L}{2} \right) \|\mathbf{v}_{j+2} - \mathbf{x}_{j+1}\|^2. \quad (168)$$

So $\alpha_x < \frac{1}{L}$ such that

$$\mathbf{v}_{j+2} = \text{prox}_{\alpha_x g}(\mathbf{x}_{j+1} - \alpha_x \nabla f(\mathbf{x}_{j+1})) \quad (169)$$

satisfies

$$F(\mathbf{v}_{j+2}) \leq c_{j+1} - \delta \|\mathbf{v}_{j+2} - \mathbf{x}_{j+1}\|^2. \quad (170)$$

Theorem 4 *Let f be a proper function with Lipschitz continuous gradients and g be proper and lower semicontinuous. Let $\Omega_1 = \{k_1, k_2, \dots, k_j, \dots\}$ and $\Omega_2 = \{m_1, m_2, \dots, m_j, \dots\}$ such that (133) in Algorithm 3 (or (144) in Algorithm 4) is executed for all $k = k_j \in \Omega_1$ and (135) in Algorithm 3 (or (149) in Algorithm 4) is executed for all $k = m_j \in \Omega_2$. For nonconvex f and nonconvex nonsmooth g , assume that (3) holds, then $\{\mathbf{x}_k\}$, $\{\mathbf{v}_k\}$ and $\{\mathbf{y}_{k_j}\}$ where $k_j \in \Omega_1$, generated by Algorithms 3 and 4, are bounded and*

1. if Ω_1 or Ω_2 is finite, then for any accumulation point $\{\mathbf{x}^*\}$ of $\{\mathbf{x}_k\}$, we have $0 \in \partial F(\mathbf{x}^*)$.
2. if Ω_1 and Ω_2 are both infinite, then for any accumulation point \mathbf{x}^* of $\{\mathbf{x}_{k_j+1}\}$, \mathbf{y}^* of $\{\mathbf{y}_{k_j}\}$ where $k_j \in \Omega_1$, and any accumulation point \mathbf{x}^* of $\{\mathbf{x}_{m_j}\}$, \mathbf{v}^* of $\{\mathbf{v}_{m_j+1}\}$ where $m_j \in \Omega_2$, we have $0 \in \partial F(\mathbf{x}^*)$, $0 \in \partial F(\mathbf{y}^*)$ and $0 \in \partial F(\mathbf{v}^*)$.

Proof From Algorithm 3 we know that if (133) (or (144) in Algorithm 4) is executed, then

$$F(\mathbf{x}_{k+1}) \leq c_k - \delta \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2, \quad (171)$$

and

$$c_{k+1} = \frac{\eta q_k c_k + F(\mathbf{x}_{k+1})}{q_{k+1}} \quad (172)$$

$$\leq \frac{\eta q_k c_k + c_k - \delta \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2}{q_{k+1}} \quad (173)$$

$$= c_k - \frac{\delta \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2}{q_{k+1}}. \quad (174)$$

If (135) (or (144) in Algorithm 4) is executed, then

$$F(\mathbf{x}_{k+1}) \leq F(\mathbf{v}_{k+1}) \leq c_k - \delta \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2, \quad (175)$$

and

$$c_{k+1} \leq c_k - \frac{\delta \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2}{q_{k+1}}. \quad (176)$$

From $F(\mathbf{x}_{k+1}) \leq c_k \leq A_k = \frac{\sum_{i=1}^k F(\mathbf{x}_i)}{k}$ we can have that $F(\mathbf{x}_{k+1})$ and c_k are bounded by induction. By assumption (3) we know that $\{\mathbf{x}_k\}$ is bounded. From $F(\mathbf{v}_{k+1}) \leq c_k$ we know \mathbf{v}_{k+1} is bounded if \mathbf{v}_{k+1} is computed.

From the definitions of Ω_1 and Ω_2 , we have

$$c_{k_j+1} \leq c_{k_j} - \frac{\delta \|\mathbf{x}_{k_j+1} - \mathbf{y}_{k_j}\|^2}{q_{k_j+1}}, k_j \in \Omega_1 \quad (177)$$

$$c_{m_j+1} \leq c_{m_j} - \frac{\delta \|\mathbf{v}_{m_j+1} - \mathbf{x}_{m_j}\|^2}{q_{m_j+1}}, m_j \in \Omega_2 \quad (178)$$

$$\Omega_1 \cup \Omega_2 = \{1, 2, 3, \dots\}, \Omega_1 \cap \Omega_2 = \emptyset. \quad (179)$$

From the definition of q_k we have

$$q_{k+1} = 1 + \sum_{i=1}^k \eta^i = \sum_{i=0}^k \eta^i \leq \sum_{i=0}^{\infty} \eta^i = \frac{1}{1-\eta}, \quad (180)$$

702 So we have

$$703 \delta(1 - \eta) \|\mathbf{x}_{k_j+1} - \mathbf{y}_{k_j}\|^2 \leq \frac{\delta \|\mathbf{x}_{k_j+1} - \mathbf{y}_{k_j}\|^2}{q_{k_j+1}} \leq c_{k_j} - c_{k_j+1}, \quad (181)$$

$$704 \delta(1 - \eta) \|\mathbf{v}_{m_j+1} - \mathbf{x}_{m_j}\|^2 \leq \frac{\delta \|\mathbf{v}_{m_j+1} - \mathbf{x}_{m_j}\|^2}{q_{m_j+1}} \leq c_{m_j} - c_{m_j+1}. \quad (182)$$

705 where $k_j \in \Omega_1, m_j \in \Omega_2$. Summing over $j = 1, \dots, \infty$, we have

$$706 \delta(1 - \eta) \sum_{j=1}^{\infty} (\|\mathbf{x}_{k_j+1} - \mathbf{y}_{k_j}\|^2 + \|\mathbf{v}_{m_j+1} - \mathbf{x}_{m_j}\|^2) \leq c_1 - F^*. \quad (183)$$

707 where $k_j \in \Omega_1, m_j \in \Omega_2$, F^* is the same function value at all the accumulation points and remark
708 that $F(\mathbf{x}_k) \leq c_k$ in Lemma 2, $\Omega_1 \cup \Omega_2 = \{1, 2, 3, \dots\}$, $\Omega_1 \cap \Omega_2 = \emptyset$ and for a fixed k , either (174)
709 or (176) holds. So we have

$$710 \sum_{j=1}^{\infty} (\|\mathbf{x}_{k_j+1} - \mathbf{y}_{k_j}\|^2 + \|\mathbf{v}_{m_j+1} - \mathbf{x}_{m_j}\|^2) \leq \frac{c_1 - F^*}{\delta(1 - \eta)} < \infty. \quad (184)$$

711 We consider three cases one by one.

712 (1) Ω_2 is finite. In this case, there exists K_0 such that (133) (or (144) in Algorithm 4) is executed
713 for all $k > K_0$. So

$$714 \sum_{k=K_0}^{\infty} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2 < \infty, \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2 \rightarrow 0. \quad (185)$$

715 From the boundness of $\{\mathbf{x}_k\}$ we have that $\{\mathbf{y}_k\}$ is bounded because $\|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2 \rightarrow 0$. Let \mathbf{y}^* be
716 any accumulation point of $\{\mathbf{y}_k\}$, say $\{\mathbf{y}_{k_l}\} \rightarrow \mathbf{y}^*$ as $l \rightarrow \infty$. From $\|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2 \rightarrow 0$ we have
717 $\{\mathbf{x}_{k_l+1}\} \rightarrow \mathbf{y}^*$ as $l \rightarrow \infty$.

718 From the optimality condition of (132) and $\mathbf{x}_{k+1} = \mathbf{z}_{k+1}$ we have

$$719 0 \in \nabla f(\mathbf{y}_{k_l}) + \frac{1}{\alpha_y} (\mathbf{x}_{k_l+1} - \mathbf{y}_{k_l}) + \partial g(\mathbf{x}_{k_l+1}) \quad (186)$$

$$720 = \nabla f(\mathbf{x}_{k_l+1}) + \nabla f(\mathbf{y}_{k_l}) - \nabla f(\mathbf{x}_{k_l+1}) + \frac{1}{\alpha_y} (\mathbf{x}_{k_l+1} - \mathbf{y}_{k_l}) + \partial g(\mathbf{x}_{k_l+1}), \quad (187)$$

721 So we have

$$722 -\nabla f(\mathbf{y}_{k_l}) + \nabla f(\mathbf{x}_{k_l+1}) - \frac{1}{\alpha_y} (\mathbf{x}_{k_l+1} - \mathbf{y}_{k_l}) \in \partial F(\mathbf{x}_{k_l+1}), \quad (188)$$

723 and

$$724 \left\| \nabla f(\mathbf{y}_{k_l}) - \nabla f(\mathbf{x}_{k_l+1}) + \frac{1}{\alpha_y} (\mathbf{x}_{k_l+1} - \mathbf{y}_{k_l}) \right\| \leq \left(\frac{1}{\alpha_y} + L \right) \|\mathbf{x}_{k_l+1} - \mathbf{y}_{k_l}\| \rightarrow 0, \quad (189)$$

725 as $l \rightarrow \infty$.

726 From (132) and $\mathbf{x}_{k+1} = \mathbf{z}_{k+1}$ we have

$$727 \langle \nabla f(\mathbf{y}_{k_l}), \mathbf{x}_{k_l+1} - \mathbf{y}_{k_l} \rangle + \frac{1}{2\alpha_y} \|\mathbf{x}_{k_l+1} - \mathbf{y}_{k_l}\|^2 + g(\mathbf{x}_{k_l+1}) \quad (190)$$

$$728 \leq \langle \nabla f(\mathbf{y}_{k_l}), \mathbf{y}^* - \mathbf{y}_{k_l} \rangle + \frac{1}{2\alpha_y} \|\mathbf{y}^* - \mathbf{y}_{k_l}\|^2 + g(\mathbf{y}^*). \quad (191)$$

729 So

$$730 \limsup_{l \rightarrow \infty} g(\mathbf{x}_{k_l+1}) \leq g(\mathbf{y}^*). \quad (192)$$

731 From the definition of lower semicontinuous of g we have

$$732 \liminf_{l \rightarrow \infty} g(\mathbf{x}_{k_l+1}) \geq g(\mathbf{y}^*). \quad (193)$$

756 So we have

$$757 \lim_{l \rightarrow \infty} g(\mathbf{x}_{k_l+1}) = g(\mathbf{y}^*). \quad (194)$$

758 Because f is continuously differentiable, we have

$$759 \lim_{l \rightarrow \infty} F(\mathbf{x}_{k_l+1}) = F(\mathbf{y}^*). \quad (195)$$

760 Similar to Theorem 1 we have

$$761 0 \in \partial F(\mathbf{y}^*). \quad (196)$$

762 From $\|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2 \rightarrow 0$ we know that $\{\mathbf{x}_k\}$ and $\{\mathbf{y}_k\}$ have the same accumulation point. So for any accumulation point \mathbf{x}^* of $\{\mathbf{x}_k\}$ we have

$$763 0 \in \partial F(\mathbf{x}^*). \quad (197)$$

764 (2) Ω_1 is finite. In this case, there exists K_0 such that (135) (or (149) in Algorithm 4) is executed for all $k > K_0$. So

$$765 \sum_{k=K_0}^{\infty} \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 < \infty, \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2 \rightarrow 0. \quad (198)$$

766 Similar to Theorem 1, for any accumulation point \mathbf{x}^* of $\{\mathbf{x}_k\}$ we have

$$767 0 \in \partial F(\mathbf{x}^*). \quad (199)$$

768 (3) Ω_1 and Ω_2 are both infinite. In this case

$$769 \|\mathbf{x}_{k_j+1} - \mathbf{y}_{k_j}\|^2 \rightarrow 0, \|\mathbf{v}_{m_j+1} - \mathbf{x}_{m_j}\|^2 \rightarrow 0. \quad (200)$$

770 where $k_j \in \Omega_1, m_j \in \Omega_2$. From the boundness of $\{\mathbf{x}_k\}$ we know \mathbf{y}_{k_j} is bounded where $k_j \in \Omega_1$. From cases 1 and 2, we know that for any accumulation point \mathbf{y}^* of $\{\mathbf{y}_{k_j}\}$, $k_j \in \Omega_1$ and any accumulation point \mathbf{x}^* of $\{\mathbf{x}_{m_j}\}$, $m_j \in \Omega_2$, we have $0 \in \partial F(\mathbf{y}^*)$ and $0 \in \partial F(\mathbf{x}^*)$. Because $\{\mathbf{x}_{k_j+1}\}$ and $\{\mathbf{y}_{k_j}\}$ have the same accumulation point for $k_j \in \Omega_1$, $\{\mathbf{v}_{m_j+1}\}$ and $\{\mathbf{x}_{m_j}\}$ have the same accumulation point for $m_j \in \Omega_2$. So for any accumulation point \mathbf{x}^* of $\{\mathbf{x}_{k_j+1}\}$, $k_j \in \Omega_1$, and any accumulation point \mathbf{v}^* of $\{\mathbf{v}_{m_j+1}\}$, $m_j \in \Omega_2$, $0 \in \partial F(\mathbf{x}^*)$, $0 \in \partial F(\mathbf{v}^*)$. ■

778 4 Numerical Results: Sparse PCA

779 Principal Component Analysis (PCA) is a basic technique for finding low-dimensional representations. But it has a drawback of lack of interpretability. Sparse PCA is a common approach to find interpretable principal components and has been applied successfully in areas such as bioinformatics [7]. One of the most popular approaches for solving Sparse PCA is the Generalized Power Method (GPower) [8]. It first solves the following problem (201), then adds a post-processing step. We focus here on the time consuming problem (201), which is an optimization problem on the Stiefel manifold:

$$780 \min_{X^T X = I} f(X) = -\frac{1}{2} \sum_{j=1}^m \sum_{i=1}^d [\mu_j |\mathbf{a}_i^T \mathbf{x}_j| - \gamma_j]_+^2, \quad (201)$$

781 where $X \in \mathbb{R}^{n \times m}$, n is the sample size, m is the desired number of PCA component, $A \in \mathbb{R}^{n \times d}$ is the data matrix, d is the sample dimension and γ controls the sparsity. $[x]_+ = \max\{x, 0\}$. We set $\mu_j = 1, \gamma_j = 0.2$ for all $1 \leq j \leq m$, and test with different m 's.

782 We compare monotone APG (mAPG) and nonmonotone APG (nmAPG) with Proximal Gradient Method (PG), GPower and the Curvilinear search method (CurviLS) [9], the state-of-art algorithm on the Stiefel manifold. The performance of PG and Inertial Forward-Backward (IFB) is similar. So we omit to list the result of IFB here. We test the performance on the breast cancer data set¹, which contains 295 samples of 8241 dimensions. All the algorithms are terminated when $\|Df(X)\|_{\infty} < 0.1$ or the number of iterations exceeds 3000, where $Df(X) := \nabla f(X) - X(\nabla f(X))^T X$ is the projected gradient onto the tangent planes. We test the machine learning performance by the sparsity

783 ¹Data available at <http://cbio.ensmp.fr/ljacob/documents/overlasso-package.tgz>

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Table 1: Comparisons of APG, PG, GPower and CurveLS on the Sparse PCA problem. The quantities include number of iterations, computing time (in seconds), sparsity (percentage of zeros) and adjusted variance. We pursuit high sparsity and variance. They are averaged over 10 runs.

m	Method	#Iter.	Time	sparsity	var
40	GPower	1557	697	0.5341	0.5532
	PG	1554	695	0.5341	0.5532
	CurviLS	647	318	0.5343	0.5541
	mAPG	275	268	0.5342	0.5536
	nmAPG	385	202	0.5341	0.5539
60	GPower	1315	711	0.5992	0.6048
	PG	1316	716	0.5992	0.6048
	CurviLS	790	474	0.5991	0.6047
	mAPG	268	322	0.5994	0.6049
	nmAPG	364	225	0.5994	0.6049
80	GPower	1574	1012	0.6457	0.6367
	PG	1575	1009	0.6457	0.6367
	CurviLS	941	662	0.6455	0.6366
	mAPG	262	371	0.6457	0.6370
	nmAPG	391	282	0.6459	0.6373

and the adjusted variance [10]. In PG and APG, we set the stepsize $\alpha = 100$. $f(X)$ in (201) is a concave function and any stepsize can ensure that $F(\mathbf{v}_{k+1}) \leq F(\mathbf{x}_k) - \delta \|\mathbf{v}_{k+1} - \mathbf{x}_k\|^2$ holds. So we choose a large stepsize to make it close to that of GPower, which can be viewed as using a stepsize of ∞ .

Table 1 shows the related result. We can note that APG-type algorithms are much faster than PG and GPower. mAPG needs fewer iterations while nmAPG needs less time. On the one hand, this indicates that the monitor-corrector step in mAPG takes effect. On the other hand, the cost of each iteration in mAPG is almost twice than that of nmAPG. This means that in nmAPG $F(\mathbf{z}_{k+1}) \leq c_k - \delta \|\mathbf{z}_{k+1} - \mathbf{y}_k\|^2$ holds almost in all iterations and accordingly \mathbf{v}_k is not computed in most of the time. We can also see that APG-type algorithms are faster than CurviLS, demonstrating that APG is a competitive method for optimization on the Stiefel manifold.

References

- [1] R.T. Rockafellar & R. Wets, Variational Analysis. Springer, 1998.
- [2] J. Bolte, S. Sabach & M. Teboulle. Proximal alternating linearized minimization for non-convex and nonsmooth problems. Mathematical Programming, 146(1-2):459-494, 2014.
- [3] H. Attouch, J. Bolte, P. Redont & A. Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: an approach based on the Kurdyka-Łojasiewicz inequality. Mathematics of Operations Research, 35:438-457, 2010.
- [4] A. Beck & M. Teboulle. Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems. IEEE Transactions on Image Processing, 18(11):2419-2434, 2009.
- [5] P. Frankel, G. Garrigos & J. Peypouquet. Splitting methods with variable metric for kurdykajosiewicz functions and general convergence rates. Journal of Optimization Theory and Applications, 165:874-900, 2014.
- [6] H. Zhang & W.W. Hager, A nonmonotone line search technique and its application to unconstrained optimization. SIAM J. Optimization, 14:1043-1056, 2004
- [7] D. Lee, W. Lee, Y. Lee & Y. Pawitan. Super-sparse principal component analyses for high-throughput genomic data. BMC Bioinformatics, 11(1):296, 2010.
- [8] M. Journee, Y. Nesterov, P. Richtarik & R. Sepulchre. Generalized power method for sparse pricipal component analysis. The Journal of Machine Learning Rearch, 11:517-553, 2010.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

[9] Z. Wen & W. Yin. A feasible method for optimization with orthogonality constraints. *Mathematical Programming*, 142:397-434, 2013.

[10] H. Zou, T. Hastie & R. Tibshirani. Sparse principal component analysis. *Journal of computational and graphical statistics*, 15(2):265-286, 2006.