

A. Overview of the Appendix

In Section B, we derive the optimal solution U^* for (2.4), which was used in Section 2. Next, we clarify the implementation details of Algorithm 1 in Section C. We illustrate more experiments in Section D. In Section F, we prove our theoretical results and some preliminary lemmas are given in Section E.

B. Optimal solution U^* for (2.4)

The optimal solution U for (2.4) is given by the first order optimality condition:

$$\frac{\partial \tilde{h}(Y, D, U)}{\partial U} = U + \lambda_3(Y^\top Y U - Y^\top D) = 0,$$

which implies

$$\begin{aligned} U^* &= \left(\frac{1}{\lambda_3} I_p + Y^\top Y \right)^{-1} Y^\top D \\ &= \lambda_3 \sum_{j=0}^{+\infty} (-\lambda_3 Y^\top Y)^j Y^\top D \\ &= \lambda_3 Y^\top \left[\sum_{j=0}^{+\infty} (-\lambda_3 Y Y^\top)^j \right] D \\ &= Y^\top \left(\frac{1}{\lambda_3} I_p + Y Y^\top \right)^{-1} D. \end{aligned}$$

Then, its transpose is given by

$$U^{*\top} = D^\top \left(\frac{1}{\lambda_3} I_p + \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^\top \right)^{-1} Y.$$

Note that \mathbf{u}_i is the column of U^\top . So for each $i \in [n]$,

$$\begin{aligned} \mathbf{u}_i^* &= D^\top \left(\frac{1}{\lambda_3} I_p + \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^\top \right)^{-1} \mathbf{y}_i \\ &= \frac{1}{n} D^\top \left(\frac{1}{\lambda_3 n} I_p + \frac{1}{n} N_n \right)^{-1} \mathbf{y}_i, \end{aligned}$$

where we denote $N_n = \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^\top$.

Also, we have

$$\begin{aligned} Y U^{*\top} &= Y \left(\frac{1}{\lambda_3} I_p + Y^\top Y \right)^{-1} Y^\top D \\ &= \lambda_3 Y \left(I_p + \lambda_3 Y^\top Y \right)^{-1} Y^\top D \\ &= \lambda_3 Y \left[\sum_{j=0}^{+\infty} (-\lambda_3 Y^\top Y)^j \right] Y^\top D \\ &= \lambda_3 \sum_{j=0}^{+\infty} (-\lambda_3)^j (Y Y^\top)^{j+1} D \\ &= D - \left(I_p + \lambda_3 Y Y^\top \right)^{-1} D \\ &= D - \left(I_p + \lambda_3 \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^\top \right)^{-1} D \\ &= D - \frac{1}{n} \left(\frac{1}{n} I_p + \frac{\lambda_3}{n} N_n \right)^{-1} D. \end{aligned}$$

Thus,

$$\begin{aligned} h(Y, D) &= \sum_{i=1}^n \frac{1}{2} \left\| \frac{1}{n} D^\top \left(\frac{1}{\lambda_3 n} I_p + \frac{1}{n} N_n \right)^{-1} \mathbf{y}_i \right\|_2^2 \\ &\quad + \frac{\lambda_3}{2} \left\| \frac{1}{n} \left(\frac{1}{n} I_p + \frac{\lambda_3}{n} N_n \right)^{-1} D \right\|_F^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \frac{1}{2} \left\| D^\top \left(\frac{1}{\lambda_3 n} I_p + \frac{1}{n} N_n \right)^{-1} \mathbf{y}_i \right\|_2^2 \\ &\quad + \frac{\lambda_3}{2n^2} \left\| \left(\frac{1}{n} I_p + \frac{\lambda_3}{n} N_n \right)^{-1} D \right\|_F^2. \end{aligned}$$

C. Algorithm Details

Algorithm 2 Solving \mathbf{v} and \mathbf{e}

Require: $D \in \mathbb{R}^{p \times d}$, $\mathbf{z} \in \mathbb{R}^p$, parameters λ_1 and λ_2

Ensure: Optimal \mathbf{v} and \mathbf{e} .

- 1: Set $\mathbf{e} = \mathbf{0}$.
- 2: **repeat**
- 3: Update \mathbf{v} :

$$\mathbf{v} = \left(D^\top D + \frac{1}{\lambda_1} I \right)^{-1} D^\top (\mathbf{z} - \mathbf{e}).$$

- 4: Update \mathbf{e} :

$$\mathbf{e} = \mathcal{S}_{\lambda_2/\lambda_1}[\mathbf{z} - D\mathbf{v}].$$

- 5: **until** convergence
-

For Algorithm 2, we set a threshold $\epsilon = 10^{-3}$. Let $\{\mathbf{v}', \mathbf{e}'\}$ and $\{\mathbf{v}'', \mathbf{e}''\}$ be the two consecutive iterates. If the maximum of $\|\mathbf{v}' - \mathbf{v}''\|_2 / \|\mathbf{v}'\|_2$ and $\|\mathbf{e}' - \mathbf{e}''\|_2 / \|\mathbf{e}'\|_2$ is less than ϵ , then we stop Algorithm 2.

Algorithm 3 Solving D

Require: $D \in \mathbb{R}^{p \times d}$ in the previous iteration, accumulation matrix M , A and B , parameters λ_1 and λ_3 .

Ensure: Optimal D (updated).

- 1: Denote $\hat{A} = \lambda_1 A + \lambda_3 I$ and $\hat{B} = \lambda_1 B + \lambda_3 M$.
- 2: **repeat**
- 3: **for** $j = 1$ to d **do**
- 4: Update the j th column of D :

$$\mathbf{d}_j \leftarrow \mathbf{d}_j - \frac{1}{\hat{A}_{jj}} \left(D \hat{\mathbf{a}}_j - \hat{\mathbf{b}}_j \right)$$

- 5: **end for**
- 6: **until** convergence

For Algorithm 3, we observe that a one-pass update on the dictionary D is enough for the final convergence of D , as we shown in the experiments. This is also observed in Mairal et al. (2010).

D. More Experiments

We also investigate the performance of subspace clustering on MNIST-7K and MNIST-10K. In this way, one can see how the computational time changes with the number of samples.

Table 4. Clustering accuracy (%) and computational time (seconds).

	OLRSC	ORPCA	LRR	LRR2	SSC
Mushrooms	85.09	65.26	58.44	56.38	54.16
	8.78	8.30	46.82	8.55	32 min
DNA	67.11	53.11	44.01	45.32	52.23
	2.58	2.09	23.67	1.65	3 min
Protein	43.30	40.22	40.31	40.00	44.27
	24.66	22.90	921.58	98.33	65 min
USPS	65.95	55.70	52.98	58.69	47.58
	33.93	27.01	257.25	71.15	50 min
MNIST-7K	58.04	55.40	54.77	54.27	45.56
	42.99	39.84	512.37	95.21	26 min
MNIST-10K	56.79	54.66	55.15	53.67	44.90
	67	56	24 min	153	84 min
MNIST-20K	57.74	54.10	55.23	54.53	43.91
	129	121	32 min	360	7 hours

E. Proof Preliminaries

Lemma 3 (Corollary of Thm. 4.1 (Bonnans & Shapiro, 1998)). Let $f : \mathbb{R}^p \times \mathbb{R}^q \rightarrow \mathbb{R}$. Suppose that for all $\mathbf{x} \in \mathbb{R}^p$ the function $f(\mathbf{x}, \cdot)$ is differentiable, and that f and $\nabla_{\mathbf{u}} f(\mathbf{x}, \mathbf{u})$ are continuous on $\mathbb{R}^p \times \mathbb{R}^q$. Let $\mathbf{v}(\mathbf{u})$ be

the optimal value function $\mathbf{v}(\mathbf{u}) = \min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x}, \mathbf{u})$, where \mathcal{C} is a compact subset of \mathbb{R}^p . Then $\mathbf{v}(\mathbf{u})$ is directionally differentiable. Furthermore, if for $\mathbf{u}_0 \in \mathbb{R}^q$, $f(\cdot, \mathbf{u}_0)$ has unique minimizer \mathbf{x}_0 then $\mathbf{v}(\mathbf{u})$ is differentiable in \mathbf{u}_0 and $\nabla_{\mathbf{u}} \mathbf{v}(\mathbf{u}_0) = \nabla_{\mathbf{u}} f(\mathbf{x}_0, \mathbf{u}_0)$.

Lemma 4 (Corollary of Donsker theorem (van der Vaart, 2000)). Let $F = \{f_{\theta} : \mathcal{X} \rightarrow \mathbb{R}, \theta \in \Theta\}$ be a set of measurable functions indexed by a bounded subset Θ of \mathbb{R}^d . Suppose that there exists a constant K such that

$$|f_{\theta_1}(x) - f_{\theta_2}(x)| \leq K \|\theta_1 - \theta_2\|_2,$$

for every θ_1 and θ_2 in Θ and x in \mathcal{X} . Then, F is P -Donsker. For any f in F , let us define $\mathbb{P}_n f$, $\mathbb{P} f$ and $\mathbb{G}_n f$ as

$$\begin{aligned} \mathbb{P}_n f &= \frac{1}{n} \sum_{i=1}^n f(X_i), \\ \mathbb{P} f &= \mathbb{E}[f(X)], \\ \mathbb{G}_n f &= \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f). \end{aligned}$$

Let us also suppose that for all f , $\mathbb{P} f^2 < \delta^2$ and $\|f\|_{\infty} < M$ and that the random elements X_1, X_2, \dots are Borel-measurable. Then, we have

$$\mathbb{E} \|\mathbb{G}\|_F = O(1),$$

where $\|\mathbb{G}\|_F = \sup_{f \in F} |\mathbb{G}_n f|$.

Lemma 5 (Sufficient condition of convergence for a stochastic process (Bottou, 1998)). Let (Ω, \mathcal{F}, P) be a measurable probability space, u_t , for $t \geq 0$, be the realization of a stochastic process and \mathcal{F}_t be the filtration by the past information at time t . Let

$$\delta_t = \begin{cases} 1 & \text{if } \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t] > 0, \\ 0 & \text{otherwise.} \end{cases}$$

If for all t , $u_t \geq 0$ and $\sum_{t=1}^{\infty} \mathbb{E}[\delta_t(u_{t+1} - u_t)] < \infty$, then u_t is a quasi-martingale and converges almost surely. Moreover,

$$\sum_{t=1}^{\infty} |\mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]| < +\infty \text{ a.s.}$$

Lemma 6 (Lemma 8 from Mairal et al. (2010)). Let a_t, b_t be two real sequences such that for all t , $a_t \geq 0, b_t \geq 0, \sum_{t=1}^{\infty} a_t = \infty, \sum_{t=1}^{\infty} a_t b_t < \infty, \exists K > 0$, such that $|b_{t+1} - b_t| < K a_t$. Then, $\lim_{t \rightarrow +\infty} b_t = 0$.

F. Proof Details

F.1. Proof of Boundedness

Proposition 7. Let $\{\mathbf{u}_t\}, \{\mathbf{v}_t\}, \{\mathbf{e}_t\}$ and $\{D_t\}$ be the optimal solutions produced by Algorithm 1. Then,

1. $\mathbf{v}_t, \mathbf{e}_t, \frac{1}{t}A_t$ and $\frac{1}{t}B_t$ are uniformly bounded.
2. M_t is uniformly bounded.
3. D_t is supported by some compact set \mathcal{D} .
4. \mathbf{u}_t is uniformly bounded.

Proof. Let us consider the optimization problem of solving \mathbf{v} and \mathbf{e} . As the trivial solution $\{\mathbf{v}'_t, \mathbf{e}'_t\} = \{\mathbf{0}, \mathbf{z}_t\}$ are feasible, we have

$$\tilde{\ell}_1(\mathbf{z}_t, D_{t-1}, \mathbf{v}'_t, \mathbf{e}'_t) = \lambda_2 \|\mathbf{z}_t\|_1.$$

Therefore, the optimal solution should satisfy:

$$\frac{\lambda_1}{2} \|\mathbf{z}_t - D_{t-1}\mathbf{v}_t - \mathbf{e}_t\|_2^2 + \frac{1}{2} \|\mathbf{v}_t\|_2^2 + \lambda_2 \|\mathbf{e}_t\|_1 \leq \lambda_2 \|\mathbf{z}_t\|_1,$$

which implies

$$\begin{aligned} \frac{1}{2} \|\mathbf{v}_t\|_2^2 &\leq \lambda_2 \|\mathbf{z}_t\|_1, \\ \lambda_2 \|\mathbf{e}_t\|_1 &\leq \lambda_2 \|\mathbf{z}_t\|_1. \end{aligned}$$

Since \mathbf{z}_t is uniformly bounded (Assumption 1), \mathbf{v}_t and \mathbf{e}_t are uniformly bounded.

To examine the uniform bound for $\frac{1}{t}A_t$ and $\frac{1}{t}B_t$, note that

$$\begin{aligned} \frac{1}{t}A_t &= \frac{1}{t} \sum_{i=1}^t \mathbf{v}_i \mathbf{v}_i^\top, \\ \frac{1}{t}B_t &= \frac{1}{t} \sum_{i=1}^t (\mathbf{z}_i - \mathbf{e}_i) \mathbf{v}_i^\top. \end{aligned}$$

Since for each i , $\mathbf{v}_i, \mathbf{e}_i$ and \mathbf{z}_i are uniformly bounded, $\frac{1}{t}A_t$ and $\frac{1}{t}B_t$ are uniformly bounded.

Now we derive the bound for M_t . All the information we have is:

1. $M_t = \sum_{i=1}^t \mathbf{y}_i \mathbf{u}_i^\top$ (definition of M_t).
2. $\mathbf{u}_t = (\|\mathbf{y}_t\|_2^2 + \frac{1}{\lambda_3})^{-1} (D_{t-1} - M_{t-1})^\top \mathbf{y}_t$ (closed form solution).
3. $D_t(\lambda_1 A_t + \lambda_3 I) = \lambda_1 B_t + \lambda_3 M_t$ (first order optimality condition for D_t).
4. $\frac{1}{t}A_t, \frac{1}{t}B_t, \frac{1}{t}N_t$ are uniformly upper bounded (Claim 1).
5. The smallest singular values of $\frac{1}{t}N_t$ and $\frac{1}{t}A_t$ are uniformly lower bounded away from zero (Assumption 2 and 3).

For simplicity, we write D_t as:

$$D_t = (\lambda_1 B_t + \lambda_3 M_t) Q_t^{-1}, \quad (\text{F.1})$$

where

$$Q_t = \lambda_1 A_t + \lambda_3 I.$$

Note that as we assume $\frac{1}{t}A_t$ is positive definite, Q_t is always invertible.

From the definition of M_t and (3.4), we know that

$$\begin{aligned} M_{t+1} - M_t &= \mathbf{y}_{t+1} \mathbf{u}_{t+1}^\top \\ &= \left(\|\mathbf{y}_{t+1}\|_2^2 + \frac{1}{\lambda_3} \right)^{-1} \mathbf{y}_{t+1} \mathbf{y}_{t+1}^\top (D_t - M_t) \\ &= P_t D_t - P_t M_t \\ &= P_t (\lambda_1 B_t + \lambda_3 M_t) Q_t^{-1} - P_t M_t, \end{aligned} \quad (\text{F.2})$$

where

$$P_t = \left(\|\mathbf{y}_{t+1}\|_2^2 + \frac{1}{\lambda_3} \right)^{-1} \mathbf{y}_{t+1} \mathbf{y}_{t+1}^\top.$$

By multiplying Q_t on both sides of (F.2), we have

$$M_{t+1} = (M_t - \lambda_1 P_t M_t A_t Q_t^{-1}) + \lambda_1 P_t B_t Q_t^{-1}. \quad (\text{F.3})$$

By applying the Taylor expansion on Q_t^{-1} , we have

$$Q_t^{-1} = (\lambda_1 A_t + \lambda_3 I_d)^{-1} = \frac{1}{\lambda_3} \sum_{i=0}^{+\infty} \left(-\frac{\lambda_1}{\lambda_3} A_t \right)^i.$$

Thus,

$$\begin{aligned} A_t Q_t^{-1} &= \frac{1}{\lambda_3} \sum_{i=0}^{+\infty} \left(-\frac{\lambda_1}{\lambda_3} \right)^i (A_t)^{i+1} \\ &= -\frac{1}{\lambda_1} \sum_{i=0}^{+\infty} \left(-\frac{\lambda_1}{\lambda_3} A_t \right)^{i+1} \\ &= -\frac{1}{\lambda_1} \left[\sum_{i=-1}^{+\infty} \left(-\frac{\lambda_1}{\lambda_3} A_t \right)^{i+1} - I_d \right] \\ &= -\frac{1}{\lambda_1} \left(I_d + \frac{\lambda_1}{\lambda_3} A_t \right)^{-1} + \frac{1}{\lambda_1} I_d. \end{aligned}$$

So M_{t+1} is given by

$$\begin{aligned} M_{t+1} &= (I_d - P_t) M_t \\ &\quad + \underbrace{P_t M_t \left(I_d + \frac{\lambda_1}{\lambda_3} A_t \right)^{-1}}_{W_t} + \lambda_1 P_t B_t Q_t^{-1}. \end{aligned} \quad (\text{F.4})$$

We first show that $P_t B_t Q_t^{-1}$ is uniformly bounded.

$$\begin{aligned} \|P_t B_t Q_t^{-1}\| &= \left\| P_t \left(\frac{1}{t} B_t \right) \left(\frac{1}{t} Q_t \right)^{-1} \right\| \\ &\leq \|P_t\| \cdot \left\| \frac{1}{t} B_t \right\| \cdot \left\| \left(\frac{1}{t} Q_t \right)^{-1} \right\|. \end{aligned}$$

Since we assume that $\{z_t\}$ are uniformly upper bounded (Assumption 1), there exists a constant α_1 , such that for all $t > 0$,

$$\|z_t\|_2 \leq \alpha_1.$$

So we have

$$\|P_{t+1}\| \leq \frac{\lambda_3 \alpha_1^2}{\lambda_3 \alpha_1^2 + 1}.$$

Next, as we have shown that $\frac{1}{t} B_t$ can be uniformly bounded, there exists a constant c_1 , such that for all $t > 0$,

$$\left\| \frac{1}{t} B_t \right\| \leq c_1.$$

And,

$$\begin{aligned} \left\| \left(\frac{1}{t} Q_t \right)^{-1} \right\| &= \frac{1}{\sigma_{\min} \left(\frac{1}{t} Q_t \right)} \\ &= \frac{1}{\sigma_{\min} \left(\frac{\lambda_1}{t} A_t + \frac{\lambda_3}{t} I_d \right)} \\ &= \frac{1}{\frac{\lambda_3}{t} + \lambda_1 \sigma_{\min} \left(\frac{1}{t} A_t \right)} \\ &\leq \frac{1}{\lambda_3 + \lambda_1 \beta_0}. \end{aligned}$$

Thus, $\lambda_1 P_t B_t Q_t^{-1}$ is uniformly bounded by a constant, say c_2 . That is,

$$\|\lambda_1 P_t B_t Q_t^{-1}\| \leq c_2. \quad (\text{F.5})$$

It follows that W_t can be bounded:

$$\begin{aligned} \|W_t\| &\leq \|P_t\| \cdot \|M_t\| \cdot \left\| \left(I_d + \frac{\lambda_1}{\lambda_3} A_t \right)^{-1} \right\| + c_2 \\ &\stackrel{\zeta_1}{\leq} \frac{\alpha_1^2}{\alpha_1^2 + \frac{1}{\lambda_3}} \cdot \frac{\lambda_3}{\lambda_3 + \lambda_1 \beta_0 t} \|M_t\| + c_2 \\ &\leq \frac{c_3}{t} \|M_t\| + c_2, \end{aligned} \quad (\text{F.6})$$

where ζ_1 is derived by utilizing the assumption that z is upper bounded by α_1 and the smallest singular value of $\frac{1}{t} A_t$ is lower bounded by β_0 . The last inequality always holds for some uniform constant c_3 .

From Assumption 2, we know that the singular values of $\frac{1}{t} \sum_{i=1}^t z_i z_i^\top$ should uniformly span the diagonal. Thus, there exists a constant τ , such that for all $i > 0$, $\frac{1}{\tau} \sum_{i=1}^{i+\tau} z_i z_i^\top$ is uniformly bounded away from zero with high probability.

Let $m_1 = \|M_1\|$. Now we pick a constant t^* , such that

$$\frac{c_3 \tau}{t^*} \left(\frac{1}{\alpha_0} + 1 \right) \leq 0.5. \quad (\text{F.7})$$

We also have a constant w^* , such that for all $t \leq t^*$,

$$\begin{aligned} \|W_t\| &\leq w^*, \\ \frac{c_3}{t} m_1 + 0.5 w^* + c_2 &\leq w^*. \end{aligned} \quad (\text{F.8})$$

Based on this, we first derive a bound for all $\|M_t\|$, with $t \leq t^*$. We know that there exists an integer k^* (which is a uniform constant), such that $k^*(\tau + 1) \leq t^* < (k^* + 1)(\tau + 1)$. Our strategy is to bound $\|M_t\|$ in each interval $[(k-1)(\tau+1), k(\tau+1)]$. We start our reasoning from the first interval $[1, \tau + 1]$.

It is easy to induce from (F.4) that for all $t > 0$,

$$M_{t+1} = \prod_{i=1}^t (I_p - P_i) M_1 + \sum_{j=1}^{t-1} \prod_{i=j+1}^t (I_p - P_i) W_j + W_t.$$

Thus,

$$\begin{aligned} &\|M_{\tau+1}\| \\ &= \left\| \prod_{i=1}^{\tau} (I_p - P_i) M_1 + \sum_{j=1}^{\tau-1} \prod_{i=j+1}^{\tau} (I_p - P_i) W_j + W_{\tau} \right\| \\ &\leq \left\| \prod_{i=1}^{\tau} (I_p - P_i) M_1 \right\| + \left\| \sum_{j=1}^{\tau-1} \prod_{i=j+1}^{\tau} (I_p - P_i) W_j + W_{\tau} \right\| \\ &\stackrel{\zeta_1}{\leq} \left\| \prod_{i=1}^{\tau} (I_p - P_i) \right\| \cdot \|M_1\| + \tau w^* \\ &\stackrel{\zeta_2}{\leq} (1 - \alpha_0) m_1 + \tau w^*. \end{aligned}$$

Here, ζ_1 holds because $\left\| \prod_{i=j+1}^{\tau} (I_p - P_i) \right\| \leq 1$ for all $j \in [\tau - 1]$. ζ_2 holds because the singular values of P_i 's have span over the diagonal so the largest singular value of $\prod_{i=1}^{\tau} (I_p - P_i)$ is $1 - \alpha_0$, where α_0 is the lower bound for all z_i 's (see Assumption 1).

For $M_{2(\tau+1)}$, we can similarly obtain

$$\|M_{2(\tau+1)}\| \leq (1 - \alpha_0)^2 m_1 + (1 - \alpha_0) \tau w^* + \tau w^*.$$

More generally, for any integer $k \leq k^*$,

$$\begin{aligned} \|M_{k(\tau+1)}\| &\leq (1-\alpha_0)^k m_1 + \sum_{j=0}^{k-1} (1-\alpha_0)^j \tau w^* \\ &\leq m_1 + \frac{\tau w^*}{\alpha_0}. \end{aligned}$$

Hence, we obtain a uniform bound for $\|M_{k(\tau+1)}\|$, with $k \in [k^*]$. For any $i \in ((k-1)(\tau+1), k(\tau+1))$, they can simply be bounded by

$$\begin{aligned} \|M_i\| &\leq m_1 + \frac{\tau w^*}{\alpha_0} + (i - (k-1)(\tau+1))w^* \\ &\leq m_1 + \frac{\tau w^*}{\alpha_0} + \tau w^*. \end{aligned}$$

Therefore, for all the current M_t 's, we can bound them as follows:

$$\|M_t\| \leq m_1 + \frac{\tau w^*}{\alpha_0} + \tau w^*, \quad \forall t = 1, 2, \dots, t^*. \quad (\text{F.9})$$

From (F.8) and (F.9), we know that for all $t \leq t^*$,

$$\begin{aligned} \|W_t\| &\leq w^*, \\ \|M_t\| &\leq m_1 + \frac{\tau w^*}{\alpha_0} + \tau w^*. \end{aligned}$$

Next, we show that the bounds still hold for $\|W_{t^*+1}\|$ and $\|M_{t^*+1}\|$, which completes our induction.

For $\|M_{t^*+1}\|$, it can simply be bounded in the same way as aforementioned because all the W_t 's are bounded by w^* for $t < t^* + 1$. That is,

$$\begin{aligned} \|M_{t^*+1}\| &\leq \|M_{k^*(\tau+1)}\| + (t^* + 1 - k^*(\tau+1))w^* \\ &\leq m_1 + \frac{\tau w^*}{\alpha_0} + \tau w^*. \end{aligned} \quad (\text{F.10})$$

For $\|W_{t^*+1}\|$, from (F.6), we know

$$\begin{aligned} \|W_{t^*+1}\| &\leq \frac{c_3}{t^*+1} \|M_{t^*+1}\| + c_2 \\ &\leq \frac{c_3}{t^*+1} (m_1 + \frac{\tau w^*}{\alpha_0} + \tau w^*) + c_2 \\ &= \frac{c_3 m_1}{t^*+1} + \frac{c_3 \tau}{t^*+1} (\frac{1}{\alpha_0} + 1) w^* + c_2 \\ &\stackrel{\zeta_1}{\leq} \frac{c_3 m_1}{t^*+1} + 0.5 w^* + c_2 \\ &\stackrel{\zeta_2}{\leq} w^*. \end{aligned} \quad (\text{F.11})$$

Here, ζ_1 is derived by utilizing (F.7) and ζ_2 is derived by (F.8).

From (F.10) and (F.11), we know that the bound for $\|M_t\|$ and $\|W_t\|$'s, with $t \leq t^*$, still holds for $t = t^* + 1$. Thus

we complete the induction and conclude that for all $t > 0$, we have

$$\begin{aligned} \|M_t\| &\leq m_1 + \frac{\tau w^*}{\alpha_0} + \tau w^*, \\ \|W_t\| &\leq w^*. \end{aligned}$$

Thus, M_t is uniformly bounded.

From (F.1), we know that

$$\begin{aligned} D_t &= \lambda_1 B_t (\lambda_1 A_t + \lambda_3 I_d)^{-1} + \lambda_3 M_t (\lambda_1 A_t + \lambda_3 I_d)^{-1} \\ &= \lambda_1 \left(\frac{1}{t} B_t \right) \left(\frac{\lambda_1}{t} A_t + \frac{\lambda_3}{t} I_d \right)^{-1} \\ &\quad + \frac{\lambda_3}{t} M_t \left(\frac{\lambda_1}{t} A_t + \frac{\lambda_3}{t} I_d \right)^{-1}. \end{aligned}$$

Since $\frac{1}{t} A_t$, $\frac{1}{t} B_t$ and M_t are all uniformly bounded, D_t is also uniformly bounded.

By examining the closed form of \mathbf{u}_t , and note that we have proved the uniform boundedness of D_t and M_t , we conclude that $\{\mathbf{u}_t\}$ are uniformly bounded. \square

Corollary 8. *Let \mathbf{v}_t , \mathbf{e}_t , \mathbf{u}_t and D_t be the optimal solutions produced by Algorithm 1.*

1. $\tilde{\ell}(\mathbf{z}_t, D_t, \mathbf{v}_t, \mathbf{e}_t)$ and $\ell(\mathbf{z}_t, D_t)$ are uniformly bounded.
2. $\frac{1}{t} \tilde{h}(Z, D, U)$ is uniformly bounded.
3. The surrogate function $g_t(D_t)$ defined in (3.5) is uniformly bounded and Lipschitz.

Proof. To show Claim 1, we just need to examine the definition of $\tilde{\ell}(\mathbf{z}_t, D_t, \mathbf{v}_t, \mathbf{e}_t)$ and notice that \mathbf{z}_t , D_t , \mathbf{v}_t and \mathbf{e}_t are all uniformly bounded. This implies that $\tilde{\ell}(\mathbf{z}_t, D_t, \mathbf{v}_t, \mathbf{e}_t)$ is uniformly bounded and so is $\ell(\mathbf{z}_t, D_t)$. Likewise, we show that $\frac{1}{t} \tilde{h}(Z, D, U)$ is uniformly bounded.

The uniform boundedness of $g_t(D_t)$ follows immediately as $\tilde{\ell}(\mathbf{z}_t, D_t, \mathbf{v}_t, \mathbf{e}_t)$ and $\frac{1}{t} \tilde{h}(Z, D, U)$ are both uniformly bounded. To show that $g_t(D)$ is Lipschitz, we show that the gradient of $g_t(D)$ is uniformly bounded for all $D \in \mathcal{D}$.

$$\begin{aligned} \|\nabla g_t(D)\|_F &= \left\| \lambda_1 D \left(\frac{A_t}{t} + \frac{\lambda_3}{t} I_d \right) - \lambda_1 \frac{B_t}{t} - \frac{\lambda_3}{t} M_t \right\|_F \\ &\leq \lambda_1 \|D\|_F \left(\left\| \frac{A_t}{t} \right\|_F + \left\| \frac{\lambda_3}{t} I_d \right\|_F \right) \\ &\quad + \lambda_1 \left\| \frac{B_t}{t} \right\|_F + \left\| \frac{\lambda_3}{t} M_t \right\|_F. \end{aligned}$$

Notice that each term on the right side of the inequality is uniformly bounded. Thus the gradient of $g_t(D)$ is uniformly bounded and $g_t(D)$ is Lipschitz. \square

Proposition 9. Let $D \in \mathcal{D}$ and denote the minimizer of $\tilde{\ell}(\mathbf{z}, D, \mathbf{v}, \mathbf{e})$ as:

$$\{\mathbf{v}', \mathbf{e}'\} = \arg \min_{\mathbf{v}, \mathbf{e}} \tilde{\ell}(\mathbf{z}, D, \mathbf{v}, \mathbf{e}).$$

Then, the function $\ell(\mathbf{z}, L)$ is continuously differentiable and

$$\nabla_D \ell(\mathbf{z}, D) = (D\mathbf{v}' + \mathbf{e}' - \mathbf{z})\mathbf{v}'^\top.$$

Furthermore, $\ell(\mathbf{z}, \cdot)$ is uniformly Lipschitz.

Proof. By fixing the variable \mathbf{z} , the function $\tilde{\ell}$ can be seen as a mapping:

$$\begin{aligned} \mathbb{R}^{d+p} \times \mathcal{D} &\rightarrow \mathbb{R} \\ ([\mathbf{v}; \mathbf{e}], D) &\mapsto \tilde{\ell}(\mathbf{z}, D, \mathbf{v}, \mathbf{e}). \end{aligned}$$

It is easy to show that for all $[\mathbf{v}; \mathbf{e}] \in \mathbb{R}^{d+p}$, $\tilde{\ell}(\mathbf{z}, \cdot, \mathbf{v}, \mathbf{e})$ is differentiable. Also $\tilde{\ell}(\mathbf{z}, \cdot, \cdot, \cdot)$ is continuous on $\mathbb{R}^{d+p} \times \mathcal{D}$. $\nabla_D \tilde{\ell}(\mathbf{z}, D, \mathbf{v}, \mathbf{e}) = (D\mathbf{v} + \mathbf{e} - \mathbf{z})\mathbf{v}^\top$ is continuous on $\mathbb{R}^{d+p} \times \mathcal{D}$. $\forall D \in \mathcal{D}$, since $\tilde{\ell}(\mathbf{z}, D, \mathbf{v}, \mathbf{e})$ is strongly convex w.r.t. \mathbf{v} , it has a unique minimizer $\{\mathbf{v}', \mathbf{e}'\}$. Thus Lemma 3 applies and we prove that $\ell(\mathbf{z}, D)$ is differentiable in D and

$$\nabla_D \ell(\mathbf{z}, D) = (D\mathbf{v}' + \mathbf{e}' - \mathbf{z})\mathbf{v}'^\top.$$

Since every term in $\nabla_D \ell(\mathbf{z}, D)$ is uniformly bounded (Assumption 1 and Proposition 7), we conclude that the gradient of $\ell(\mathbf{z}, D)$ is uniformly bounded, implying that $\ell(\mathbf{z}, D)$ is uniformly Lipschitz w.r.t. D . \square

Corollary 10. Let $f_t(D)$ be the empirical loss function defined in (2.6). Then $f_t(D)$ is uniformly bounded and Lipschitz.

Proof. Since $\ell(\mathbf{z}, L)$ can be uniformly bounded (Corollary 8), we only need to show that $\frac{1}{t}h(Z, D)$ is uniformly bounded. Note that we have derived the form for $h(Z, D)$ as follows:

$$\begin{aligned} \frac{1}{t}h(Z, D) &= \frac{1}{t^3} \sum_{i=1}^t \frac{1}{2} \left\| D^\top \left(\frac{1}{\lambda_3 t} I_p + \frac{1}{t} N_t \right)^{-1} \mathbf{z}_i \right\|_2^2 \\ &\quad + \frac{\lambda_3}{2t^3} \left\| \left(\frac{1}{t} I_p + \frac{\lambda_3}{t} N_t \right)^{-1} D \right\|_F^2 \end{aligned}$$

where $N_t = \sum_{i=1}^t \mathbf{z}_i \mathbf{z}_i^\top$. Since every term in the above equation can be uniformly bounded, $h(Z, D)$ is uniformly bounded and so is $f_t(D)$.

To show that $f_t(D)$ is uniformly Lipschitz, we show that its gradient can be uniformly bounded.

$$\begin{aligned} &\nabla f_t(D) \\ &= \frac{1}{t} \sum_{i=1}^t \nabla \ell(\mathbf{z}_i, D) + \frac{1}{t} \nabla h(Z, D) \\ &= \frac{1}{t} \sum_{i=1}^t (D\mathbf{v}_i + \mathbf{e}_i - \mathbf{z}_i)\mathbf{v}_i^\top \\ &\quad + \frac{1}{t^3} \sum_{i=1}^t \left(\frac{1}{\lambda_3 t} I_p + \frac{1}{t} N_t \right)^{-1} \mathbf{z}_i \mathbf{z}_i^\top \left(\frac{1}{\lambda_3 t} I_p + \frac{1}{t} N_t \right)^{-1} D \\ &\quad + \frac{\lambda_3}{t^3} \left(\frac{1}{t} I_p + \frac{\lambda_3}{t} N_t \right)^{-2} D. \end{aligned}$$

Then the Frobenius norm of $\nabla f_t(D)$ can be bounded by:

$$\begin{aligned} &\|\nabla f_t(D)\|_F \\ &\leq \frac{1}{t} \sum_{i=1}^t \|D\mathbf{v}_i + \mathbf{e}_i - \mathbf{z}_i\|_2 \cdot \|\mathbf{v}_i\|_2 \\ &\quad + \frac{1}{t^3} \sum_{i=1}^t \left\| \left(\frac{1}{\lambda_3 t} I_p + \frac{1}{t} N_t \right)^{-1} \right\|_F^2 \cdot \|\mathbf{z}_i\|_2^2 \cdot \|D\|_F \\ &\quad + \frac{\lambda_3}{t^3} \left\| \left(\frac{1}{t} I_p + \frac{\lambda_3}{t} N_t \right)^{-1} \right\|_F^2 \cdot \|D\|_F. \end{aligned}$$

One can easily check that the right side of the inequality is uniformly bounded. Thus $\|\nabla f_t(D)\|_F$ is uniformly bounded, implying that $f_t(D)$ is uniformly Lipschitz. \square

F.2. Proof of P-Donsker

Proposition 11. Let $f'_t(D) = \frac{1}{t} \sum_{i=1}^t \ell(\mathbf{z}_i, D)$ and $f(D)$ be the expected loss function defined in (2.8). Then we have

$$\mathbb{E}[\sqrt{t} \|f'_t - f\|_\infty] = \mathcal{O}(1).$$

Proof. Let us consider $\{\ell(\mathbf{z}, D)\}$ as a set of measurable functions indexed by $D \in \mathcal{D}$. As we showed in Proposition 7, \mathcal{D} is a compact set. Also, we have proved that $\ell(\mathbf{z}, D)$ is uniformly Lipschitz over D (Proposition 9). Thus, $\{\ell(\mathbf{z}, D)\}$ is P-Donsker (see the definition in Lemma 4). Furthermore, as $\ell(\mathbf{z}, D)$ is non-negative and uniformly bounded, so is $\ell^2(\mathbf{z}, D)$. So we have $\mathbb{E}_z[\ell^2(\mathbf{z}, D)]$ being uniformly bounded. Since we have verified all the hypotheses in Lemma 4, we obtain the result that

$$\mathbb{E}[\sqrt{t} \|f'_t - f\|_\infty] = \mathcal{O}(1).$$

\square

F.3. Proof of convergence of $g_t(D)$

Theorem 12 (Convergence of the surrogate function $g_t(D_t)$). *The surrogate function $g_t(D_t)$ we defined in (3.5) converges almost surely, where D_t is the solution produced by Algorithm 1.*

Proof. Note that $g_t(D_t)$ can be viewed as a stochastic positive process since every term in $g_t(D_t)$ is non-negative and the samples are drawn randomly. We define

$$u_t = g_t(D_t).$$

To show the convergence of u_t , we need to bound the difference of u_{t+1} and u_t :

$$\begin{aligned} & u_{t+1} - u_t \\ &= g_{t+1}(D_{t+1}) - g_t(D_t) \\ &= g_{t+1}(D_{t+1}) - g_{t+1}(D_t) + g_{t+1}(D_t) - g_t(D_t) \\ &= g_{t+1}(D_{t+1}) - g_{t+1}(D_t) \\ &\quad + \frac{1}{t+1} \ell(\mathbf{z}_{t+1}, D_t) - \frac{1}{t+1} g'_t(D_t) \\ &\quad + \left[\frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{\lambda_3}{2(t+1)} \|D_t - M_{t+1}\|_F^2 \right. \\ &\quad \left. - \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 - \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \right] \\ &= g_{t+1}(D_{t+1}) - g_{t+1}(D_t) + \frac{f'_t(D_t) - g'_t(D_t)}{t+1} \\ &\quad + \frac{\ell(\mathbf{z}_{t+1}, D_t) - f'_t(D_t)}{t+1} \\ &\quad + \left[\frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{\lambda_3}{2(t+1)} \|D_t - M_{t+1}\|_F^2 \right. \\ &\quad \left. - \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 - \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \right]. \end{aligned} \quad (\text{F.12})$$

Here,

$$g'_t(D_t) = \frac{1}{t} \sum_{i=1}^t \tilde{\ell}(\mathbf{z}_i, D, \mathbf{v}_i, \mathbf{e}_i). \quad (\text{F.13})$$

First, we bound the last four terms. We have

$$\begin{aligned} & \frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 - \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 \\ &= \frac{-1}{t(t+1)} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2 \\ &\leq \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2. \end{aligned} \quad (\text{F.14})$$

And

$$\begin{aligned} & \frac{\lambda_3}{2(t+1)} \|D_t - M_{t+1}\|_F^2 - \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \\ &= \frac{-\lambda_3}{2t(t+1)} \|D_t - M_t\|_F^2 + \frac{\lambda_3}{2(t+1)} \|\mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top\|_F^2 \\ &\quad - \frac{\lambda_3}{t+1} \text{Tr}((D_t - M_t)^\top \mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top) \\ &= \frac{-\lambda_3}{2t(t+1)} \|D_t - M_t\|_F^2 + \frac{\lambda_3}{2(t+1)} \|\mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top\|_F^2 \\ &\quad - \frac{\lambda_3}{t+1} \left(\|\mathbf{z}_{t+1}\|_2^2 + \frac{1}{\lambda_3} \right) \|\mathbf{u}_{t+1}\|_2^2 \\ &\leq \frac{1}{t+1} \left(\frac{\lambda_3}{2} \|\mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top\|_F^2 - (\lambda_3 \|\mathbf{z}_{t+1}\|_2^2 + 1) \|\mathbf{u}_{t+1}\|_2^2 \right) \\ &\leq \frac{1}{t+1} \left(-\frac{\lambda_3}{2} \|\mathbf{z}_{t+1}\|_2^2 \|\mathbf{u}_{t+1}\|_2^2 - \|\mathbf{u}_{t+1}\|_2^2 \right), \end{aligned} \quad (\text{F.15})$$

where the first equality is derived by the fact that $M_{t+1} = M_t + \mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top$, and the second equality is derived by the closed form solution of \mathbf{u}_{t+1} (see (3.4)).

Combining (F.14) and (F.15), we know that

$$\begin{aligned} & \frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 - \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 \\ &\quad + \frac{\lambda_3}{2(t+1)} \|D_t - M_{t+1}\|_F^2 - \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \\ &\leq \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2 + \frac{1}{t+1} \left(-\frac{\lambda_3}{2} \|\mathbf{z}_{t+1}\|_2^2 \|\mathbf{u}_{t+1}\|_2^2 \right. \\ &\quad \left. - \|\mathbf{u}_{t+1}\|_2^2 \right) \\ &= \frac{1}{t+1} \left(-\frac{\lambda_3}{2} \|\mathbf{z}_{t+1}\|_2^2 \|\mathbf{u}_{t+1}\|_2^2 - \frac{1}{2} \|\mathbf{u}_{t+1}\|_2^2 \right) \leq 0. \end{aligned}$$

Therefore,

$$\begin{aligned} u_{t+1} - u_t &\leq g_{t+1}(D_{t+1}) - g_{t+1}(D_t) + \frac{1}{t+1} \ell(\mathbf{z}_{t+1}, D_t) \\ &\quad - \frac{1}{t+1} g'_t(D_t) \\ &= g_{t+1}(D_{t+1}) - g_{t+1}(D_t) + \frac{f'_t(D_t) - g'_t(D_t)}{t+1} \\ &\quad + \frac{\ell(\mathbf{z}_{t+1}, D_t) - f'_t(D_t)}{t+1} \\ &\leq \frac{\ell(\mathbf{z}_{t+1}, D_t) - f'_t(D_t)}{t+1}, \end{aligned}$$

where $f'_t(D)$ is defined in Proposition 11, and the last inequality holds because D_{t+1} is the minimizer of $g_{t+1}(D)$ and $g'_t(D)$ is a surrogate function of $f'_t(D)$.

Let \mathcal{F}_t be the filtration of the past information. We take the

expectation on the above equation conditioned on \mathcal{F}_t :

$$\begin{aligned}\mathbb{E}[u_{t+1} - u_t \mid \mathcal{F}_t] &\leq \frac{\mathbb{E}[\ell(\mathbf{z}_{t+1}, D_t) \mid \mathcal{F}_t] - f'_t(D_t)}{t+1} \\ &\leq \frac{f(D_t) - f'_t(D_t)}{t+1} \\ &\leq \frac{\|f - f'_t\|_\infty}{t+1}.\end{aligned}$$

From Proposition 11, we know

$$\mathbb{E}[\|f - f'_t\|_\infty] = \mathcal{O}\left(\frac{1}{\sqrt{t}}\right).$$

Thus,

$$\begin{aligned}\mathbb{E}[\mathbb{E}[u_{t+1} - u_t \mid \mathcal{F}_t]^+] &= \mathbb{E}[\max\{\mathbb{E}[u_{t+1} - u_t \mid \mathcal{F}_t], 0\}] \\ &\leq \frac{c}{\sqrt{t}(t+1)},\end{aligned}$$

where c is some constant.

Now let us define the index set

$$\mathcal{T} = \{t \mid \mathbb{E}[u_{t+1} - u_t \mid \mathcal{F}_t] > 0\},$$

and the indicator

$$\delta_t = \begin{cases} 1, & \text{if } t \in \mathcal{T}, \\ 0, & \text{otherwise.} \end{cases}$$

We have

$$\begin{aligned}\sum_{t=1}^{\infty} \mathbb{E}[\delta_t(u_{t+1} - u_t)] &= \sum_{t \in \mathcal{T}} \mathbb{E}[u_{t+1} - u_t] \\ &= \sum_{t \in \mathcal{T}} \mathbb{E}[\mathbb{E}[u_{t+1} - u_t \mid \mathcal{F}_t]] \\ &= \sum_{t=1}^{\infty} \mathbb{E}[\mathbb{E}[u_{t+1} - u_t \mid \mathcal{F}_t]^+] \\ &\leq +\infty.\end{aligned}$$

Thus, Lemma 5 applies. That is, $g_t(D_t)$ is a quasi-martingale and converges almost surely. Moreover,

$$\sum_{t=1}^{\infty} |\mathbb{E}[u_{t+1} - u_t \mid \mathcal{F}_t]| < +\infty, \text{ a.s.} \quad (\text{F.16})$$

□

F.4. Proof of Convergence of D_t

Proposition 13. *Let $\{D_t\}_{t=1}^{\infty}$ be the basis sequence produced by the Algorithm 1. Then,*

$$\|D_{t+1} - D_t\|_F = \mathcal{O}\left(\frac{1}{t}\right). \quad (\text{F.17})$$

Proof. According the strong convexity of $g_t(D)$ (Assumption 3), we have,

$$g_t(D_{t+1}) - g_t(D_t) \geq \frac{\beta_0}{2} \|D_{t+1} - D_t\|_F^2, \quad (\text{F.18})$$

On the other hand,

$$\begin{aligned}&g_t(D_{t+1}) - g_t(D_t) \\ &= g_t(D_{t+1}) - g_{t+1}(D_{t+1}) + g_{t+1}(D_{t+1}) - g_{t+1}(D_t) \\ &\quad + g_{t+1}(D_t) - g_t(D_t) \\ &\leq g_t(D_{t+1}) - g_{t+1}(D_{t+1}) + g_{t+1}(D_t) - g_t(D_t) \\ &\stackrel{\text{def}}{=} G_t(D_{t+1}) - G_t(D_t).\end{aligned} \quad (\text{F.19})$$

Note that the inequality is derived by the fact that $g_{t+1}(D_{t+1}) - g_{t+1}(D_t) \leq 0$, as D_{t+1} is the minimizer of $g_{t+1}(D)$. We denote $g_t(D) - g_{t+1}(D)$ by $G_t(D)$.

By a simple calculation, we obtain the gradient of $G_t(D)$:

$$\begin{aligned}\nabla G_t(D) &= \nabla g_t(D) - \nabla g_{t+1}(D) \\ &= \frac{1}{t} \left[D(\lambda_1 A_t + \lambda_3 I_d) - (\lambda_1 B_t + \lambda_3 M_t) \right] \\ &\quad - \frac{1}{t+1} \left[D(\lambda_1 A_{t+1} + \lambda_3 I_d) - (\lambda_1 B_{t+1} + \lambda_3 M_{t+1}) \right] \\ &= \frac{1}{t} \left[D \left(\lambda_1 A_t + \lambda_3 I_d - \frac{\lambda_1 t}{t+1} A_{t+1} - \frac{\lambda_3 t}{t+1} I_d \right) \right. \\ &\quad \left. + \frac{\lambda_1 t}{t+1} B_{t+1} - \lambda_1 B_t + \frac{\lambda_3 t}{t+1} M_{t+1} - \lambda_3 M_t \right] \\ &= \frac{1}{t} \left[D \left(\frac{\lambda_1}{t+1} A_{t+1} - \lambda_1 \mathbf{v}_{t+1} \mathbf{v}_{t+1}^\top + \frac{\lambda_3}{t+1} I_d \right) \right. \\ &\quad \left. + \lambda_1 (\mathbf{z}_{t+1} - \mathbf{e}_{t+1}) \mathbf{v}_{t+1}^\top - \frac{\lambda_1}{t+1} B_{t+1} \right. \\ &\quad \left. + \lambda_3 \mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top - \frac{\lambda_3}{t+1} M_{t+1} \right]\end{aligned}$$

So the Frobenius norm of $\nabla G_t(D)$ follows immediately:

$$\begin{aligned}
 & \|\nabla G_t(D)\|_F \\
 & \leq \frac{1}{t} \left[\|D\|_F \left(\lambda_1 \left\| \frac{A_{t+1}}{t+1} \right\|_F + \lambda_1 \|\mathbf{v}_{t+1} \mathbf{v}_{t+1}^\top\|_F \right. \right. \\
 & \quad \left. \left. + \frac{\lambda_3}{t+1} \|I_d\|_F \right) + \lambda_1 \|(\mathbf{z}_{t+1} - \mathbf{e}_{t+1}) \mathbf{v}_{t+1}^\top\|_F \right. \\
 & \quad \left. + \lambda_1 \left\| \frac{B_{t+1}}{t+1} \right\|_F + \lambda_3 \|\mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top\|_F \right. \\
 & \quad \left. + \frac{\lambda_3}{t+1} \|M_{t+1}\|_F \right] \\
 & = \frac{1}{t} \left[\|D\|_F \left(\lambda_1 \left\| \frac{A_{t+1}}{t+1} \right\|_F + \lambda_1 \|\mathbf{v}_{t+1} \mathbf{v}_{t+1}^\top\|_F \right) \right. \\
 & \quad \left. + \lambda_1 \|(\mathbf{z}_{t+1} - \mathbf{e}_{t+1}) \mathbf{v}_{t+1}^\top\|_F \right. \\
 & \quad \left. + \lambda_1 \left\| \frac{B_{t+1}}{t+1} \right\|_F + \lambda_3 \|\mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top\|_F \right] \\
 & \quad + \frac{\lambda_3}{t(t+1)} \left[\|I_d\|_F + \|M_{t+1}\|_F \right].
 \end{aligned}$$

We know from Proposition 7 that all the terms in the above equation are uniformly bounded. Thus, there exist constants c_1 , c_2 and c_3 , such that

$$\|\nabla G_t(D)\|_F \leq \frac{1}{t} [c_1 \|D\|_F + c_2] + \frac{c_3}{t(t+1)}.$$

According to the first order Taylor expansion,

$$\begin{aligned}
 & G_t(D_{t+1}) - G_t(D_t) \\
 & = \text{Tr} \left((D_{t+1} - D_t)^\top \nabla G_t(\alpha D_t + (1-\alpha) D_{t+1}) \right) \\
 & \leq \|D_{t+1} - D_t\|_F \cdot \|\nabla G_t(\alpha D_t + (1-\alpha) D_{t+1})\|_F,
 \end{aligned}$$

where α is a constant between 0 and 1. According to Proposition 7, D_t and D_{t+1} are uniformly bounded, so $\alpha D_t + (1-\alpha) D_{t+1}$ is uniformly bounded. Thus, there exists a constant c_4 , such that

$$\|\nabla G_t(\alpha D_t + (1-\alpha) D_{t+1})\|_F \leq \frac{c_4}{t} + \frac{c_3}{t(t+1)},$$

resulting in

$$G_t(D_{t+1}) - G_t(D_t) \leq \left(\frac{c_4}{t} + \frac{c_3}{t(t+1)} \right) \|D_{t+1} - D_t\|_F.$$

Combining (F.18), (F.19) and the above equation, we have

$$\|D_{t+1} - D_t\|_F = \frac{2c_4}{\beta_0} \cdot \frac{1}{t} + \frac{2c_3}{\beta_0} \cdot \frac{1}{t(t+1)}.$$

F.5. Proof for convergence of $f_t(D_t)$

Theorem 14 (Convergence of $f_t(D_t)$). *Let $f_t(D_t)$ be the empirical loss function defined in (2.6) and D_t be the solution produced by the Algorithm 1. Let $b_t = g_t(D_t) - f_t(D_t)$. Then, b_t converges almost surely to 0. Thus, $f_t(D_t)$ converges almost surely to the same limit of $g_t(D_t)$.*

Proof. Let $f'_t(D)$ and $g'_t(D)$ be those defined in Proposition 11 and Theorem 12 respectively. Then,

$$\begin{aligned}
 b_t & = g_t(D_t) - f_t(D_t) \\
 & = g'_t(D_t) - f'_t(D_t) + \left[\frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \right. \\
 & \quad \left. - \frac{1}{t^3} \sum_{i=1}^t \frac{1}{2} \left\| D_t^\top \left(\frac{1}{\lambda_3 t} I_p + \frac{1}{t} N_t \right)^{-1} \mathbf{z}_i \right\|_2^2 \right. \\
 & \quad \left. - \frac{\lambda_3}{2t^3} \left\| \left(\frac{1}{t} I_p + \frac{\lambda_3}{t} N_t \right)^{-1} D_t \right\|_F^2 \right] \\
 & = g'_t(D_t) - f'_t(D_t) + q_t(D_t),
 \end{aligned}$$

where $q_t(D_t)$ denotes the last four terms. Combining F.12, we have

$$\begin{aligned}
 \frac{b_t}{t+1} & = \frac{g'_t(D_t) - f'_t(D_t)}{t+1} + \frac{q_t(D_t)}{t+1} \\
 & = g_{t+1}(D_{t+1}) - g_{t+1}(D_t) + \frac{\ell(\mathbf{z}_{t+1}, D_t) - f'_t(D_t)}{t+1} \\
 & \quad + u_t - u_{t+1} \\
 & \quad + \left[\frac{q_t(D_t)}{t+1} + \frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 \right. \\
 & \quad \left. + \frac{\lambda_3}{2(t+1)} \|D_t - M_{t+1}\|_F^2 - \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 \right. \\
 & \quad \left. - \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \right].
 \end{aligned}$$

Note that we naturally have

$$\begin{aligned}
 \frac{q_t(D_t)}{t+1} & \leq \frac{1}{t(t+1)} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{\lambda_3}{2t(t+1)} \|D_t - M_t\|_F^2 \\
 & \leq \frac{1}{t(t+1)} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{c}{2t(t+1)},
 \end{aligned}$$

where the second inequality holds as D_t and M_t are both uniformly bounded (see Proposition 7). \square

Also, from (F.14), we know

$$\begin{aligned} & \frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 - \frac{1}{t} \sum_{i=1}^t \|\mathbf{u}_i\|_2^2 \\ &= \frac{-1}{t(t+1)} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2. \end{aligned}$$

And from (F.15)

$$\begin{aligned} & \frac{\lambda_3}{2(t+1)} \|D_t - M_{t+1}\|_F^2 - \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \\ & \leq \frac{1}{t+1} \left(-\frac{\lambda_3}{2} \|\mathbf{z}_{t+1}\|_2^2 \|\mathbf{u}_{t+1}\|_2^2 - \|\mathbf{u}_{t+1}\|_2^2 \right). \end{aligned}$$

Thus,

$$\begin{aligned} & \frac{g_t(D_t)}{t+1} + \frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 \\ & + \frac{\lambda_3}{2(t+1)} \|D_t - M_{t+1}\|_F^2 \\ & - \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 - \frac{\lambda_3}{2t} \|D_t - M_t\|_F^2 \\ & \leq \frac{c}{2t(t+1)} + \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2 \\ & + \frac{1}{t+1} \left(-\frac{\lambda_3}{2} \|\mathbf{z}_{t+1}\|_2^2 \|\mathbf{u}_{t+1}\|_2^2 - \|\mathbf{u}_{t+1}\|_2^2 \right) \\ & = \frac{c}{2t(t+1)} - \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2 \\ & - \frac{\lambda_3}{2(t+1)} \|\mathbf{z}_{t+1}\|_2^2 \|\mathbf{u}_{t+1}\|_2^2 \\ & \leq \frac{c}{2t(t+1)}. \end{aligned}$$

Therefore,

$$\begin{aligned} & \frac{b_t}{t+1} \\ & \leq g_{t+1}(D_{t+1}) - g_{t+1}(D_t) + \frac{\ell(\mathbf{z}_{t+1}, D_t) - f'_t(D_t)}{t+1} \\ & + u_t - u_{t+1} + \frac{c}{2t(t+1)} \\ & \leq \frac{\ell(\mathbf{z}_{t+1}, D_t) - f'_t(D_t)}{t+1} + u_t - u_{t+1} + \frac{c}{2t(t+1)}. \end{aligned}$$

By taking the expectation conditioned on the past information \mathcal{F}_t , we have

$$\begin{aligned} \frac{b_t}{t+1} & \leq \frac{f(D_t) - f_t(D_t)}{t+1} + \mathbb{E}[u_t - u_{t+1} | \mathcal{F}_t] + \frac{c}{2t(t+1)} \\ & \leq \frac{c_1}{\sqrt{t}(t+1)} + |\mathbb{E}[u_t - u_{t+1} | \mathcal{F}_t]| + \frac{c}{2t(t+1)}, \end{aligned}$$

where the second inequality holds by applying Proposition 11. Thus,

$$\begin{aligned} & \sum_{t=1}^{\infty} \frac{b_t}{t+1} \\ & \leq \sum_{t=1}^{\infty} \frac{c_1}{\sqrt{t}(t+1)} + \sum_{t=1}^{\infty} |\mathbb{E}[u_t - u_{t+1} | \mathcal{F}_t]| \\ & + \sum_{t=1}^{\infty} \frac{c}{2t(t+1)} \\ & < +\infty. \end{aligned}$$

Here, the last inequality is derived by applying (F.16).

Next, we examine the difference between b_{t+1} and b_t :

$$\begin{aligned} & |b_{t+1} - b_t| \\ & = |g_{t+1}(D_{t+1}) - f_{t+1}(D_{t+1}) - g_t(D_t) + f_t(D_t)| \\ & \leq |g_{t+1}(D_{t+1}) - g_t(D_{t+1})| + |g_t(D_{t+1}) - g_t(D_t)| \\ & + |f_{t+1}(D_{t+1}) - f_t(D_{t+1})| + |f_t(D_{t+1}) - f_t(D_t)|. \end{aligned} \tag{F.20}$$

For the first term on the right hand side,

$$\begin{aligned} & |g_{t+1}(D_{t+1}) - g_t(D_{t+1})| \\ & = \left| g'_{t+1}(D_{t+1}) - g'_t(D_{t+1}) + \frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{u}_i\|_2^2 \right. \\ & \quad \left. - \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 + \frac{\lambda_3}{2(t+1)} \|D_{t+1} - M_{t+1}\|_F^2 \right. \\ & \quad \left. - \frac{\lambda_3}{2t} \|D_{t+1} - M_t\|_F^2 \right| \\ & = \left| g'_{t+1}(D_{t+1}) - g'_t(D_{t+1}) - \frac{1}{t(t+1)} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 \right. \\ & \quad \left. - \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2 \right. \\ & \quad \left. - \frac{\lambda_3}{2t(t+1)} \|D_{t+1} - M_t\|_F^2 - \frac{\lambda_3}{2(t+1)} \|\mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top\|_F^2 \right| \\ & \leq |g'_{t+1}(D_{t+1}) - g'_t(D_{t+1})| + \frac{1}{t(t+1)} \sum_{i=1}^t \frac{1}{2} \|\mathbf{u}_i\|_2^2 \\ & + \frac{1}{2(t+1)} \|\mathbf{u}_{t+1}\|_2^2 + \frac{\lambda_3}{2t(t+1)} \|D_{t+1} - M_t\|_F^2 \\ & + \frac{\lambda_3}{2(t+1)} \|\mathbf{z}_{t+1} \mathbf{u}_{t+1}^\top\|_F^2 \\ & \stackrel{\zeta_1}{\leq} |g'_{t+1}(D_{t+1}) - g'_t(D_{t+1})| + \frac{c_1}{t+1} \\ & = \left| \frac{1}{t+1} \ell(\mathbf{z}_{t+1}, D_{t+1}) - \frac{1}{t+1} g'_t(D_{t+1}) \right| + \frac{c_1}{t+1} \\ & \stackrel{\zeta_2}{\leq} \frac{c_2}{t+1}, \end{aligned}$$

where c_1 and c_2 are some uniform constants. Note that ζ_1 holds because all the \mathbf{u}_i 's, D_{t+1} , M_t and \mathbf{z}_{t+1} are uniformly bounded (see Proposition 7), and ζ_2 holds because $\ell(\mathbf{z}_{t+1}, D_{t+1})$ and $g'_t(D_{t+1})$ are uniformly bounded (see Corollary 8).

For the third term on the right hand side of (F.20), we can similarly derive

$$\begin{aligned} & |f_{t+1}(D_{t+1}) - f_t(D_{t+1})| \\ & \leq |f'_{t+1}(D_{t+1}) - f'_t(D_{t+1})| + \frac{c_3}{t+1} \\ & = \left| \frac{1}{t+1} \ell(\mathbf{z}_{t+1}, D_{t+1}) - \frac{1}{t+1} f'_t(D_{t+1}) \right| + \frac{c_3}{t+1} \\ & \leq \frac{c_4}{t+1}, \end{aligned}$$

where c_3 and c_4 are some uniform constants, and ζ_3 holds as $\ell(\mathbf{z}_{t+1}, D_{t+1})$ and $f'_t(D_{t+1})$ are both uniformly bounded (see Corollary 10).

From Corollary 8 and Corollary 10, we know that both $g_t(D)$ and $f_t(D)$ are uniformly Lipschitz. That is, there exists uniform constants κ_1, κ_2 , such that

$$\begin{aligned} |g_t(D_{t+1}) - g_t(D_t)| & \leq \kappa_1 \|D_{t+1} - D_t\|_F \stackrel{\zeta_4}{\leq} \frac{\kappa_3}{t+1}, \\ |f_t(D_{t+1}) - f_t(D_t)| & \leq \kappa_2 \|D_{t+1} - D_t\|_F \stackrel{\zeta_5}{\leq} \frac{\kappa_4}{t+1}. \end{aligned}$$

Here, ζ_4 and ζ_5 are derived by applying Proposition 13 and κ_3 and κ_4 are some uniform constants.

Finally, we have a bound for (F.20):

$$|b_{t+1} - b_t| \leq \frac{\kappa_0}{t+1},$$

where κ_0 is some uniform constant.

By applying Lemma 6, we conclude that $\{b_t\}$ converges to zero. That is,

$$\lim_{t \rightarrow +\infty} g_t(D_t) - f_t(D_t) = 0.$$

Since we have proved in Theorem 12 that $g_t(D_t)$ converges almost surely, we conclude that $f_t(D_t)$ converges almost surely to the same limit of $g_t(D_t)$. \square

Theorem 15 (Convergence of $f(D_t)$). *Let $f(D)$ be the expected loss function we defined in (2.8) and let D_t be the optimal solution produced by Algorithm 1. Then $f(D_t)$ converges almost surely to the same limit of $f_t(D_t)$ (or, $g_t(D_t)$).*

Proof. According to the central limit theorem, we know that $\sqrt{t}(f(D_t) - f_t(D_t))$ is bounded, implying

$$\lim_{t \rightarrow +\infty} f(D_t) - f_t(D_t) = 0, \quad a.s.$$

\square

F.6. Proof of gradient of $f(D)$

Proposition 16 (Gradient of $f(D)$). *Let $f(D)$ be the expected loss function which is defined in (2.8). Then, $f(D)$ is continuously differentiable and $\nabla f(D) = \mathbb{E}_{\mathbf{z}}[\nabla_D \ell(\mathbf{z}, D)]$. Moreover, $\nabla f(D)$ is uniformly Lipschitz on \mathcal{D} .*

Proof. We have shown in Proposition 9 that $\ell(\mathbf{z}, D)$ is continuously differentiable, $f(D)$ is also continuously differentiable and we have $\nabla f(D) = \mathbb{E}_{\mathbf{z}}[\nabla_D \ell(\mathbf{z}, D)]$.

Next, we prove the Lipschitz of $\nabla f(D)$. Let $\mathbf{v}'(\mathbf{z}', D')$ and $\mathbf{e}'(\mathbf{z}', D')$ be the minimizer of $\ell(\mathbf{z}', D', \mathbf{v}, \mathbf{e})$. Since $\tilde{\ell}(\mathbf{z}, D, \mathbf{v}, \mathbf{e})$ has a unique minimum and is continuous in $\mathbf{z}, D, \mathbf{v}$ and \mathbf{e} , $\mathbf{v}'(\mathbf{z}', D')$ and $\mathbf{e}'(\mathbf{z}', D')$ is continuous in \mathbf{z} and D .

Let $\Lambda = \{j \mid e'_j \neq 0\}$. According the first order optimality condition, we know that

$$\frac{\partial \tilde{\ell}(\mathbf{z}, D, \mathbf{v}, \mathbf{e})}{\partial \mathbf{e}} = 0,$$

which implies

$$\lambda_1(\mathbf{z} - D\mathbf{v} - \mathbf{e}) \in \lambda_2 \|\mathbf{e}\|_1.$$

Hence,

$$|(z - D\mathbf{v} - \mathbf{e})_j| = \frac{\lambda_2}{\lambda_1}, \quad \forall j \in \Lambda.$$

Since $\mathbf{z} - D\mathbf{v} - \mathbf{e}$ is continuous in \mathbf{z} and D , there exists an open neighborhood \mathcal{V} , such that for all $(\mathbf{z}'', D'') \in \mathcal{V}$, if $j \notin \Lambda$, then $|(z'' - D''\mathbf{v}'' - \mathbf{e}'')_j| < \frac{\lambda_2}{\lambda_1}$ and $e''_j = 0$. That is, the support set of \mathbf{e}' will not change.

Let us denote $H = [D \ I_p]$, $\mathbf{r} = [\mathbf{v}^\top \ \mathbf{e}^\top]^\top$ and define the function

$$\begin{aligned} \tilde{\ell}(\mathbf{z}, H, \mathbf{r}_\Lambda) & = \frac{\lambda_1}{2} \|\mathbf{z} - H_\Lambda \mathbf{r}_\Lambda\|_2^2 + \frac{1}{2} \|[I \ 0] \mathbf{r}_\Lambda\|_2^2 \\ & \quad + \lambda_2 \|[0 \ I] \mathbf{r}_\Lambda\|_1 \end{aligned}$$

Since $\tilde{\ell}(\mathbf{z}, D_\Lambda, \cdot)$ is strongly convex, there exists a uniform constant κ_1 , such that for all \mathbf{r}'_Λ ,

$$\begin{aligned} & \tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}''_\Lambda) - \tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}'_\Lambda) \\ & \geq \kappa_1 \|\mathbf{r}''_\Lambda - \mathbf{r}'_\Lambda\|_2^2 \\ & = \kappa_1 \left(\|\mathbf{v}'' - \mathbf{v}'\|_2^2 + \|\mathbf{e}''_\Lambda - \mathbf{e}'_\Lambda\|_2^2 \right). \end{aligned} \quad (\text{F.21})$$

On the other hand,

$$\begin{aligned} & \tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}'_\Lambda) - \tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}'_\Lambda) \\ & = \tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}'_\Lambda) - \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}''_\Lambda) \\ & \quad + \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}''_\Lambda) - \tilde{\ell}(\mathbf{z}'', D'_\Lambda, \mathbf{r}'_\Lambda) \\ & \leq \tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}'_\Lambda) - \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}''_\Lambda) \\ & \quad + \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}'_\Lambda) - \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}''_\Lambda), \end{aligned} \quad (\text{F.22})$$

where the last inequality holds because \mathbf{r}'' is the minimizer of $\tilde{\ell}(\mathbf{z}'', H'', \mathbf{r})$.

We shall prove that $\tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}_\Lambda) - \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}_\Lambda)$ is Lipschitz w.r.t. \mathbf{r} , which implies the Lipschitz of $\mathbf{v}'(\mathbf{z}, D)$ and $\mathbf{e}'(\mathbf{z}, D)$.

$$\begin{aligned} & \nabla_{\mathbf{r}} \left(\tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}_\Lambda) - \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}_\Lambda) \right) \\ &= \lambda_1 \left[H'_\Lambda{}^\top (H'_\Lambda - H''_\Lambda) + (H'_\Lambda - H''_\Lambda)^\top H''_\Lambda \right. \\ & \quad \left. + H'_\Lambda{}^\top (\mathbf{z}'' - \mathbf{z}') + (H''_\Lambda - H'_\Lambda)^\top \mathbf{z}'' \right]. \end{aligned}$$

Note that $\|H'_\Lambda\|_F$, $\|H''_\Lambda\|_F$ and \mathbf{z}'' are all uniformly bounded. Hence, there exists uniform constants c_1 and c_2 , such that

$$\begin{aligned} & \left\| \nabla_{\mathbf{r}} \left(\tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}_\Lambda) - \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}_\Lambda) \right) \right\|_2 \\ & \leq c_1 \|H'_\Lambda - H''_\Lambda\|_F + c_2 \|\mathbf{z}' - \mathbf{z}''\|_2, \end{aligned}$$

which implies that $\tilde{\ell}(\mathbf{z}', H'_\Lambda, \mathbf{r}_\Lambda) - \tilde{\ell}(\mathbf{z}'', H''_\Lambda, \mathbf{r}_\Lambda)$ is Lipschitz with Lipschitz constant $c(H'_\Lambda, H''_\Lambda, \mathbf{z}', \mathbf{z}'') = c_1 \|H'_\Lambda - H''_\Lambda\|_F + c_2 \|\mathbf{z}' - \mathbf{z}''\|_2$. Combining this fact with (F.21) and (F.22), we obtain

$$\kappa_1 \|\mathbf{r}''_\Lambda - \mathbf{r}'_\Lambda\|_2^2 \leq c(H'_\Lambda, H''_\Lambda, \mathbf{z}', \mathbf{z}'') \|\mathbf{r}''_\Lambda - \mathbf{r}'_\Lambda\|_2.$$

Therefore, $\mathbf{r}(\mathbf{z}, D)$ is Lipschitz and so are $\mathbf{v}(\mathbf{z}, D)$ and $\mathbf{e}(\mathbf{z}, D)$. Note that according to Proposition 9,

$$\begin{aligned} & \nabla f(D') - \nabla f(D'') \\ &= \mathbb{E}_{\mathbf{z}} \left[(H' \mathbf{r}' - \mathbf{z}) \mathbf{v}'^\top - (H'' \mathbf{r}'' - \mathbf{z}) \mathbf{v}''^\top \right] \\ &= \mathbb{E}_{\mathbf{z}} \left[H' \mathbf{r}' (\mathbf{v}' - \mathbf{v}'')^\top + (H' - H'') \mathbf{r}' \mathbf{v}''^\top \right. \\ & \quad \left. + H'' (\mathbf{r}' - \mathbf{r}'') \mathbf{v}''^\top + \mathbf{z} (\mathbf{v}'' - \mathbf{v}')^\top \right]. \end{aligned}$$

Thus,

$$\begin{aligned} & \|\nabla f(D') - \nabla f(D'')\|_F \\ & \stackrel{\zeta_1}{\leq} \mathbb{E}_{\mathbf{z}} \left[\|H' \mathbf{r}'\|_2 \|\mathbf{v}' - \mathbf{v}''\|_2 + \|H' - H''\|_F \|\mathbf{r}' \mathbf{v}''^\top\|_F \right. \\ & \quad \left. + \|H''\|_F \|\mathbf{r}' - \mathbf{r}''\|_2 \|\mathbf{v}''\|_2 + \|\mathbf{z}\|_2 \|\mathbf{v}' - \mathbf{v}''\|_2 \right] \\ & \stackrel{\zeta_2}{\leq} \mathbb{E}_{\mathbf{z}} \left[(\gamma_1 + \gamma_2 \|\mathbf{z}\|_2) \|H' - H''\|_F \right] \\ & \stackrel{\zeta_3}{\leq} \gamma_0 \|D' - D''\|_F, \end{aligned}$$

where γ_0 , γ_1 and γ_2 are all uniform constants. Here, ζ_1 holds because for any function $s(\mathbf{z})$, we have $\|\mathbb{E}_{\mathbf{z}}[s(\mathbf{z})]\|_F \leq \mathbb{E}_{\mathbf{z}}[\|s(\mathbf{z})\|_F]$. ζ_2 is derived by using the result that $\mathbf{r}(\mathbf{z}, H)$ and $\mathbf{v}(\mathbf{z}, H)$ are both Lipschitz and H' , H'' , \mathbf{r}' , \mathbf{r}'' , \mathbf{v}' and \mathbf{v}'' are all uniformly bounded. ζ_3 holds because \mathbf{z} is uniformly bounded and actually $\|H' - H''\|_F = \|D' - D''\|_F$. Thus, we complete the proof. \square

E.7. Proof of stationary point

Theorem 17 (Convergence of D_t). *Let $\{D_t\}$ be the optimal basis produced by Algorithm 1 and let $f(D)$ be the expected loss function defined in (2.8). Then D_t converges to a stationary point of $f(D)$ when t goes to infinity.*

Proof. Since $\frac{1}{t}A_t$ and $\frac{1}{t}B_t$ are uniformly bounded (Proposition 7), there exist sub-sequences of $\{\frac{1}{t}A_t\}$ and $\{\frac{1}{t}B_t\}$ that converge to A_∞ and B_∞ respectively. Then D_t will converge to D_∞ . Let W be an arbitrary matrix in $\mathbb{R}^{p \times d}$ and $\{h_k\}$ be any positive sequence that converges to zero.

As g_t is a surrogate function of f_t , for all t and k , we have

$$g_t(D_t + h_k W) \geq f_t(D_t + h_k W).$$

Let t tend to infinity, and note that $f(D) = \lim_{t \rightarrow \infty} f_t(D)$, we have

$$g_\infty(D_\infty + h_k W) \geq f(D_\infty + h_k W).$$

Note that the Lipschitz of ∇f indicates that the second derivative of $f(D)$ is uniformly bounded. By a simple calculation, we can also show that it also holds for $g_t(D)$. This fact implies that we can take the first order Taylor expansion for both $g_t(D)$ and $f(D)$ even when t tends to infinity (because the second order derivatives of them always exist). That is,

$$\begin{aligned} & \text{Tr}(h_k W^\top \nabla g_\infty(D_\infty)) + o(h_k W) \\ & \geq \text{Tr}(h_k W^\top \nabla f(D_\infty)) + o(h_k W) \end{aligned}$$

By multiplying $\frac{1}{h_k \|W\|_F}$ on both sides and note that $\{h_k\}$ is a positive sequence, it follows that

$$\begin{aligned} & \text{Tr} \left(\frac{1}{\|W\|_F} W^\top \nabla g_\infty(D_\infty) \right) + \frac{o(h_k W)}{h_k \|W\|_F} \\ & \geq \text{Tr} \left(\frac{1}{\|W\|_F} W^\top \nabla f(D_\infty) \right) + \frac{o(h_k W)}{h_k \|W\|_F}. \end{aligned}$$

Now let k go to infinity,

$$\text{Tr} \left(\frac{1}{\|W\|_F} W^\top \nabla g_\infty(D_\infty) \right) \geq \text{Tr} \left(\frac{1}{\|W\|_F} W^\top \nabla f(D_\infty) \right).$$

Note that this inequality holds for any matrix $W \in \mathbb{R}^{p \times d}$, so we actually have

$$\nabla g_\infty(D_\infty) = \nabla f(D_\infty).$$

As D_∞ is the minimizer of $g_\infty(D)$, we have

$$\nabla f(D_\infty) = \nabla g_\infty(D_\infty) = 0.$$

\square

References

- Avron, Haim, Kale, Satyen, Kasiviswanathan, Shiva Prasad, and Sindhvani, Vikas. Efficient and practical stochastic subgradient descent for nuclear norm regularization. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- Bertsekas, Dimitri P. *Nonlinear programming*. Athena Scientific, 1999.
- Bonnans, J. Frédéric and Shapiro, Alexander. Optimization problems with perturbations: A guided tour. *SIAM Review*, 40(2):228–264, 1998.
- Bottou, Léon. Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9), 1998.
- Burer, Samuel and Monteiro, Renato D. C. Local minima and convergence in low-rank semidefinite programming. *Mathematical Programming*, 103(3):427–444, 2005.
- Cai, Jian-Feng, Candès, Emmanuel J., and Shen, Zuowei. A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization*, 20(4):1956–1982, 2010.
- Candès, Emmanuel J. and Recht, Benjamin. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- Candès, Emmanuel J., Li, Xiaodong, Ma, Yi, and Wright, John. Robust principal component analysis? *Journal of the ACM*, 58(3):11, 2011.
- Defazio, Aaron, Bach, Francis R., and Lacoste-Julien, Simon. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Neural Information Processing Systems*, pp. 1646–1654, 2014.
- Elhamifar, Ehsan and Vidal, René. Sparse subspace clustering. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 2790–2797, 2009.
- Eriksson, Brian, Balzano, Laura, and Nowak, Robert D. High-rank matrix completion and subspace clustering with missing data. *CoRR*, abs/1112.5629, 2011.
- Fazel, Maryam, Hindi, Haitham, and Boyd, Stephen P. A rank minimization heuristic with application to minimum order system approximation. In *Proceedings of the American Control Conference*, volume 6, pp. 4734–4739. IEEE, 2001.
- Feng, Jiashi, Xu, Huan, and Yan, Shuicheng. Online robust PCA via stochastic optimization. In *Proceedings of the 27th Annual Conference on Neural Information Processing Systems*, pp. 404–412, 2013.
- Guo, Han, Qiu, Chenlu, and Vaswani, Namrata. An online algorithm for separating sparse and low-dimensional signal sequences from their sum. *IEEE Trans. Signal Processing*, 62(16):4284–4297, 2014.
- Hsieh, Cho-Jui and Olsen, Peder A. Nuclear norm minimization via active subspace selection. In *Proceedings of the 31st International Conference on Machine Learning*, pp. 575–583, 2014.
- Jaggi, Martin and Sulovský, Marek. A simple algorithm for nuclear norm regularized problems. In *Proceedings of the 27th International Conference on Machine Learning*, pp. 471–478, 2010.
- Lin, Zhouchen, Chen, Minming, and Ma, Yi. The augmented lagrange multiplier method for exact recovery of corrupted low-rank matrices. *arXiv preprint arXiv:1009.5055*, 2010.
- Liu, Guangcan and Li, Ping. Recovery of coherent data via low-rank dictionary pursuit. In *Proceedings of the 28th Annual Conference on Neural Information Processing Systems*, pp. 1206–1214, 2014.
- Liu, Guangcan, Lin, Zhouchen, Yan, Shuicheng, Sun, Ju, Yu, Yong, and Ma, Yi. Robust recovery of subspace structures by low-rank representation. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 35(1):171–184, 2013.
- Mairal, Julien, Bach, Francis R., Ponce, Jean, and Sapiro, Guillermo. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11:19–60, 2010.
- Mardani, Morteza, Mateos, Gonzalo, and Giannakis, Georgios B. Subspace learning and imputation for streaming big data matrices and tensors. *IEEE Trans. Signal Processing*, 63(10):2663–2677, 2015.
- Ng, Andrew Y., Jordan, Michael I., and Weiss, Yair. On spectral clustering: Analysis and an algorithm. In *Proceedings of the 15th Annual Conference on Neural Information Processing Systems*, pp. 849–856, 2001.
- Qiu, Chenlu, Vaswani, Namrata, Lois, Brian, and Hogben, Leslie. Recursive robust PCA or recursive sparse recovery in large but structured noise. *IEEE Trans. Information Theory*, 60(8):5007–5039, 2014.
- Recht, Benjamin, Fazel, Maryam, and Parrilo, Pablo A. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.

- Richtárik, Peter and Takác, Martin. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2014.
- Shen, Jie and Li, Ping. Learning structured low-rank representation via matrix factorization. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, pp. 500–509, 2016.
- Soltanolkotabi, Mahdi and Candès, Emmanuel J. A geometric analysis of subspace clustering with outliers. *The Annals of Statistics*, 40:2195–2238, 2012.
- Soltanolkotabi, Mahdi, Elhamifar, Ehsan, and Candès, Emmanuel J. Robust subspace clustering. *The Annals of Statistics*, 42(2):669–699, 2014.
- van der Vaart, A.W. *Asymptotic statistics*. Cambridge University Press, 2000.
- Vidal, René. A tutorial on subspace clustering. *IEEE Signal Processing Magazine*, 28(2):52–68, 2010.
- Wang, Yu-Xiang, Xu, Huan, and Leng, Chenlei. Provable subspace clustering: When LRR meets SSC. In *Proceedings of 27th Annual Conference on Neural Information Processing Systems*, pp. 64–72, 2013.
- Xiao, Lin and Zhang, Tong. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- Xu, Huan, Caramanis, Constantine, and Mannor, Shie. Principal component analysis with contaminated data: The high dimensional case. In *Proceedings of the 23rd Conference on Learning Theory*, pp. 490–502, 2010.