

# Scalable Robust Matrix Recovery: Frank-Wolfe Meets Proximal Methods

Cun Mu, Yuqian Zhang, John Wright, Donald Goldfarb (2017)

2020.6.19

王佳文

- 1.文章摘要
- 2.问题背景
- 3.问题描述
- 4.问题求解
- 5.数值实验

# 1. 文章摘要

从压缩和严重破坏的观测数据中恢复矩阵是稳健统计中的一个基本问题，在计算机视觉和机器学习中有丰富的应用。理论上，在一定的条件下，这个问题可以通过一个自然的凸松弛在多项式时间内求解，称为压缩主成分追踪(CPCP)。现有的许多可证明收敛的CPCP算法每次都存在**超线性**的迭代代价，这严重限制了它们对大规模问题的适用性。在本文中提出了可证明的收敛性、可伸缩性和有效的方法来求解具有每次迭代代价为**线性**的CPCP。主要结合了**FW算法**和**临近法**。在每次迭代中，利用**FW法**以及**秩一SVD来更新低秩项**，并利用**临近梯度法更新稀疏项**。讨论了收敛结果和实现细节。通过对可视化数据的数值实验，证明了该方法的实用性和可扩展性。

## 2.问题背景



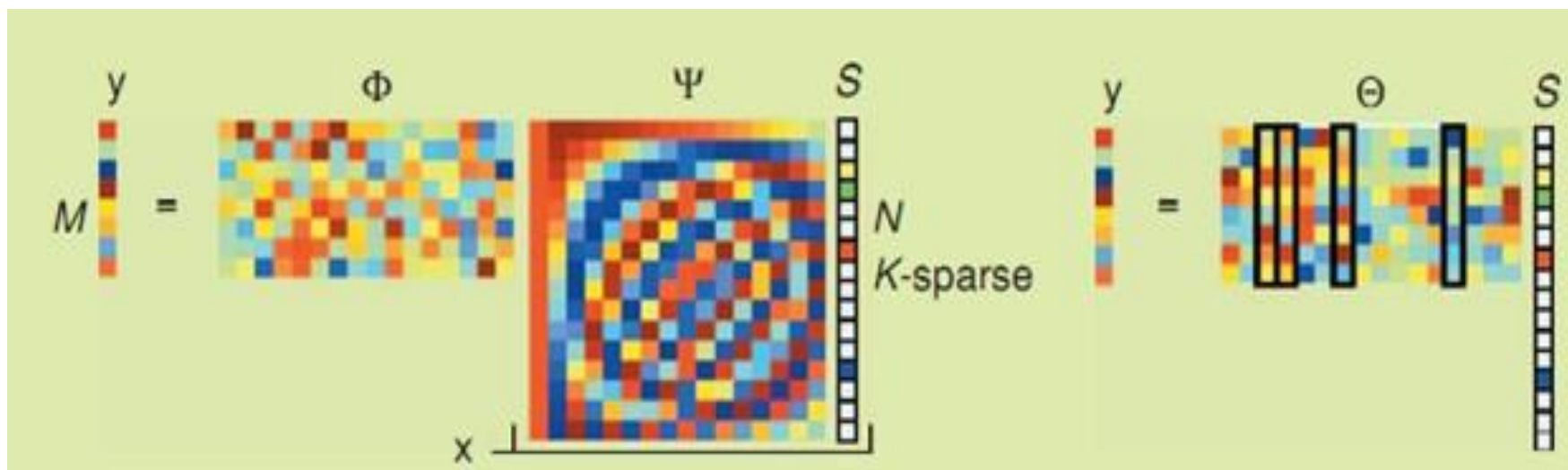
背景建模的最简单情形是从固定摄像机拍摄的视频中分离背景和前景。将视频图像序列的每一帧图像像素值拉成一个列向量，那么多个帧也就是多个列向量就组成了一个观测矩阵。由于背景比较稳定，图像序列帧与帧之间具有极大的相似性，所以仅由背景像素组成的矩阵具有低秩特性；同时由于前景是移动的物体，占据像素比例较低，故前景像素组成的矩阵具有稀疏特性。视频观测矩阵就是这两种特性矩阵的叠加，因此，可以说视频背景建模实现的过程就是低秩矩阵恢复的过程。

# 3. 问题描述

## 3.1 压缩感知简介

- 设  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ , 稀疏度为  $k$  (即含有  $k$  个非零值),  $\boldsymbol{\phi} \in \mathbb{R}^{m \times n} (m < n)$ ,  $\mathbf{y} = \boldsymbol{\phi}\mathbf{x}$  为长度为  $m$  的一维测量值。
- 压缩感知问题就是已知测量值  $\mathbf{y}$  和测量矩阵  $\boldsymbol{\phi}$  的基础上, 求解欠定方程组  $\mathbf{y} = \boldsymbol{\phi}\mathbf{x}$  得到原信号  $\mathbf{x}$ 。
- $\boldsymbol{\phi}$  的每一行可以看作是一个传感器, 它与信号相乘, 拾取了信号的一部分信息。而这一部分信息足以代表原信号, 并能找到一个算法来高概率恢复原信号。

- 一般的自然信号 $\mathbf{x}$ 本身并不是稀疏的，需要在某种稀疏基上进行稀疏表示， $\mathbf{x} = \boldsymbol{\psi}\mathbf{s}$ ， $\boldsymbol{\psi}$ 为稀疏基矩阵， $\mathbf{s}$ 为稀疏系数( $\mathbf{s}$ 只有 $k$ 个是非零值( $k \ll m$ )。压缩感知方程为 $\mathbf{y} = \boldsymbol{\phi}\mathbf{x} = \boldsymbol{\phi}\boldsymbol{\psi}\mathbf{s} = \boldsymbol{\theta}\mathbf{s}$ 。将原来的测量矩阵 $\boldsymbol{\phi}$ 变换为 $\boldsymbol{\theta} = \boldsymbol{\phi}\boldsymbol{\psi}$  (称之为传感矩阵)，解出 $\mathbf{s}$ 的逼近值 $\mathbf{s}'$ ，则原信号 $\mathbf{x}' = \boldsymbol{\psi}\mathbf{s}'$ 。



## 3.2. 问题

- 假设  $\mathbf{M}_0 \in \mathbb{R}^{m \times n}$ ,  $\mathbf{M}_0 = \mathbf{L}_0 + \mathbf{S}_0 + \mathbf{N}_0$ ,  $\mathbf{L}_0$  是一个低秩矩阵,  $\mathbf{S}_0$  是一个稀疏矩阵,  $\mathbf{N}_0$  是一个密度噪声矩阵。

- 定义:

$$(1.1) \quad \mathbf{b} = \mathbf{A}[\mathbf{M}_0] = (\langle \mathbf{A}_1, \mathbf{M}_0 \rangle, \langle \mathbf{A}_2, \mathbf{M}_0 \rangle, \dots, \langle \mathbf{A}_p, \mathbf{M}_0 \rangle)^T \in \mathbb{R}^p$$

其中  $\mathbf{A}: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  是感知算子,  $\mathbf{A}_k$  表示感知矩阵, 且:

$$\langle \mathbf{A}_k, \mathbf{M}_0 \rangle = \text{Tr}(\mathbf{M}_0^T \mathbf{A}_k)$$

- 在给定  $\mathbf{A}$  时, 我们可以有效的从  $\mathbf{b}$  中恢复  $\mathbf{L}_0$  和  $\mathbf{S}_0$  吗?

### 3.3 压缩主成分追踪(CPCP)问题

- 自然的，可以用如下方式进行优化：

$$(1.2) \quad \min_{\mathbf{L}, \mathbf{S}} \frac{1}{2} \|\mathbf{b} - \mathcal{A}[\mathbf{L} + \mathbf{S}]\|_2^2 + \lambda_L \text{rank}(\mathbf{L}) + \lambda_S \|\mathbf{S}\|_0.$$

缺点：非凸，无法直接处理

转化为：

$$(1.3) \quad \min_{\mathbf{L}, \mathbf{S}} \frac{1}{2} \|\mathbf{b} - \mathcal{A}[\mathbf{L} + \mathbf{S}]\|_2^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1$$

其中 $\|\mathbf{L}\|_*$ 定义为 $\mathbf{L}$ 的奇异值之和， $\|\mathbf{S}\|_1 = \sum_{i=1} \sum_{j=1} |\mathbf{s}_{i,j}|$

(1.3)式有时被称为压缩主成分追踪(CPCP)问题.



- 等价的, 由于:

$$\{ \mathbf{M} \in \mathbb{R}^{m \times n} \mid \mathbf{b} = \mathcal{A}[\mathbf{M}] \} = \{ \mathbf{M} \in \mathbb{R}^{m \times n} \mid \mathcal{P}_Q[\mathbf{M}] = \mathcal{P}_Q[\mathbf{M}_0] \}$$

其中  $Q \subseteq \mathbb{R}^{m \times n}$  是由感知矩阵  $\{\mathbf{A}_i\}_{i=1}^p$  张成的子空间,  $\mathcal{P}_Q$  表示在子空间上的投影算子。

- 可将原问题转化为:

$$(1.4) \quad \min_{\mathbf{L}, \mathbf{S}} f(\mathbf{L}, \mathbf{S}) \doteq \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}_0]\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1$$

注意问题的目标函数不可微且可行集无界。

Recently, CPCP and its close variants have been studied for different sensing operators  $\mathcal{A}$  (or equivalently different subspaces  $\mathcal{Q}$ ). In specific, [2, 3, 4, 5, 6] consider the case where a subset  $\Omega \subseteq \{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$  of the entries of  $\mathbf{M}_0$  is observed. Then CPCP can be reduced to

$$(1.5) \quad \min_{\mathbf{L}, \mathbf{S}} \quad \frac{1}{2} \|\mathcal{P}_\Omega[\mathbf{L} + \mathbf{S} - \mathbf{M}_0]\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1,$$

where  $\mathcal{P}_\Omega[\cdot]$  denotes the orthogonal projection onto the linear space of matrices supported on  $\Omega$ , i.e.,  $\mathcal{P}_\Omega[\mathbf{M}_0](i, j) = (\mathbf{M}_0)_{ij}$  if  $(i, j) \in \Omega$  and  $\mathcal{P}_\Omega[\mathbf{M}_0](i, j) = 0$  otherwise. [1] studies the case where each  $\mathcal{A}_k$  is an i.i.d.  $\mathcal{N}(0, 1)$  matrix, which is equivalent (in distribution) to saying that we choose a linear subspace  $\mathcal{Q}$  uniformly at random from the set of all  $p$ -dimensional subspaces of  $\mathbb{R}^{m \times n}$  and observe  $\mathcal{P}_\mathcal{Q}[\mathbf{M}_0]$ . Accordingly, all the above provide theoretical guarantees for CPCP, under fairly mild conditions, to produce accurate estimates of  $\mathbf{L}_0$  and  $\mathcal{P}_\Omega[\mathbf{S}_0]$  (or  $\mathbf{S}_0$ ), even when the number of measurements  $p$  is substantially less than  $mn$ .

## 4. 问题求解

### 4.1 FW算法

- 设 $h$ 是一紧集上的可微凸函数.

$$(2.1) \quad \text{minimize } h(\boldsymbol{x}) \quad \text{subject to } \boldsymbol{x} \in \mathcal{D} \subseteq \mathbb{R}^n$$

假设 $\nabla h$ 是L利普希茨,

在第k次迭代, 在点  $\boldsymbol{x}^k$  处线性化目标函数 $h$ 得:

$$(2.2) \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathcal{D}, \quad \|\nabla h(\boldsymbol{x}) - \nabla h(\boldsymbol{y})\| \leq L \|\boldsymbol{x} - \boldsymbol{y}\|$$

$$(2.3) \quad h(\boldsymbol{v}) \approx h(\boldsymbol{x}^k) + \langle \nabla h(\boldsymbol{x}^k), \boldsymbol{v} - \boldsymbol{x}^k \rangle$$

$$(2.4) \quad \boldsymbol{v}^k \in \arg \min_{\boldsymbol{v} \in \mathcal{D}} \langle \nabla h(\boldsymbol{x}^k), \boldsymbol{v} \rangle$$

然后在可行下降方向  $\mathbf{v}^k - \mathbf{x}^k$  上采取:

$$(2.5) \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \frac{2}{k+2}(\mathbf{v}^k - \mathbf{x}^k).$$

---

**Algorithm 1** Frank-Wolfe method for problem (2.1)

---

```
1: Initialization:  $\mathbf{x}^0 \in \mathcal{D}$ ;  
2: for  $k = 0, 1, 2, \dots$  do  
3:    $\mathbf{v}^k \in \operatorname{argmin}_{\mathbf{v} \in \mathcal{D}} \langle \mathbf{v}, \nabla h(\mathbf{x}^k) \rangle$ ;  
4:    $\gamma = \frac{2}{k+2}$ ;  
5:    $\mathbf{x}^{k+1} = \mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k)$ ;  
6: end for
```

---

---

**Algorithm 2** Frank-Wolfe method for problem (2.1) with general updating scheme

---

```
1: Initialization:  $\mathbf{x}^0 \in \mathcal{D}$ ;  
2: for  $k = 0, 1, 2, \dots$  do  
3:    $\mathbf{v}^k \in \operatorname{argmin}_{\mathbf{v} \in \mathcal{D}} \langle \mathbf{v}, \nabla h(\mathbf{x}^k) \rangle$ ;  
4:    $\gamma = \frac{2}{k+2}$  ;  
5:   Update  $\mathbf{x}^{k+1}$  to some point in  $\mathcal{D}$  such that  $h(\mathbf{x}^{k+1}) \leq h(\mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k))$ ;  
6: end for
```

---

THEOREM 2.1. Let  $\mathbf{x}^*$  be an optimal solution to (2.1). For  $\{\mathbf{x}^k\}$  generated by Algorithm 2, we have for  $k = 0, 1, 2, \dots$ ,

$$(2.9) \quad h(\mathbf{x}^k) - h(\mathbf{x}^*) \leq \frac{2LD^2}{k+2}.$$

*Proof.* For  $k = 0, 1, 2, \dots$ , we have

$$(2.10) \quad \begin{aligned} h(\mathbf{x}^{k+1}) &\leq h(\mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k)) \\ &\leq h(\mathbf{x}^k) + \gamma \langle \nabla h(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + \frac{L\gamma^2}{2} \|\mathbf{v}^k - \mathbf{x}^k\|^2 \\ &\leq h(\mathbf{x}^k) + \gamma \langle \nabla h(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + \frac{\gamma^2 LD^2}{2} \end{aligned}$$

$$(2.11) \quad \begin{aligned} &\leq h(\mathbf{x}^k) + \gamma \langle \nabla h(\mathbf{x}^k), \mathbf{x}^* - \mathbf{x}^k \rangle + \frac{\gamma^2 LD^2}{2} \\ &\leq h(\mathbf{x}^k) + \gamma(h(\mathbf{x}^*) - h(\mathbf{x}^k)) + \frac{\gamma^2 LD^2}{2}, \end{aligned}$$

其中  $D = \max_{\mathbf{x}, \mathbf{y} \in \mathcal{D}} \|\mathbf{x} - \mathbf{y}\|$  表示可行集  $\mathcal{D}$  的直径

Rearranging terms in (2.11), one obtains that for  $k = 0, 1, 2, \dots$ ,

$$(2.12) \quad h(\mathbf{x}^{k+1}) - h(\mathbf{x}^\star) \leq (1 - \gamma) (h(\mathbf{x}^k) - h(\mathbf{x}^\star)) + \frac{\gamma^2 L D^2}{2}.$$

Therefore, by mathematical induction, it can be verified that

$$h(\mathbf{x}^k) - h(\mathbf{x}^\star) \leq \frac{2LD^2}{k+2}, \quad \text{for } k = 1, 2, 3, \dots$$

REMARK 1. *Note that the constant in the rate of convergence depends on the Lipschitz constant  $L$  of  $h$  and the diameter  $\mathcal{D}$ .*

## 4.2 对偶间隙

对偶间隙:  $d(\mathbf{x}^k) = \langle \mathbf{x}^k - \mathbf{v}^k, \nabla h(\mathbf{x}^k) \rangle$

- 在  $h(\mathbf{x}^k) - h(\mathbf{x}^*)$  上提供一个有用的上界

$$\begin{aligned} h(\mathbf{x}^k) - h(\mathbf{x}^*) &\leq -\langle \mathbf{x}^* - \mathbf{x}^k, \nabla h(\mathbf{x}^k) \rangle \\ &\leq -\min_{\mathbf{v}} \langle \mathbf{v} - \mathbf{x}^k, \nabla h(\mathbf{x}^k) \rangle = \langle \mathbf{x}^k - \mathbf{v}^k, \nabla h(\mathbf{x}^k) \rangle = d(\mathbf{x}^k) \end{aligned}$$

THEOREM 2.2. *Let  $\{\mathbf{x}^k\}$  be the sequence generated by Algorithm 2. Then for any  $K \geq 1$ , there exists  $1 \leq \tilde{k} \leq K$  such that*

$$(2.15) \quad d(\mathbf{x}^{\tilde{k}}) \leq \frac{6LD^2}{K+2}.$$

REMARK 2. *The convergence rate for the duality gap matches the one for  $h(\mathbf{x}^k) - h(\mathbf{x}^*)$  (see (2.9)), which suggests that the upper bound  $d(\mathbf{x}^k)$  can serve as a practical stopping criterion.*



## 4.3 几个重要的方法

*Minimizing a linear function over the nuclear norm ball.* Since the dual norm of the nuclear norm is the operator norm, i.e.,  $\|\mathbf{Y}\| = \max_{\|\mathbf{X}\|_* \leq 1} \langle \mathbf{Y}, \mathbf{X} \rangle$ , the optimization problem

$$(2.19) \quad \text{minimize}_{\mathbf{X}} \quad \langle \mathbf{Y}, \mathbf{X} \rangle \quad \text{subject to } \|\mathbf{X}\|_* \leq 1$$

has optimal value  $-\|\mathbf{Y}\|$ . One minimizer is the rank-one matrix  $\mathbf{X}^* = -\mathbf{u}\mathbf{v}^\top$ , where  $\mathbf{u}$  and  $\mathbf{v}$  are the left- and right- singular vectors corresponding to the leading singular value of  $\mathbf{Y}$ , and can be efficiently computed (e.g. using power method).

*Minimizing a linear function over the  $\ell_1$  ball.* Since the dual norm of the  $\ell_1$  norm is the  $\ell_\infty$  norm, i.e.,  $\|\mathbf{Y}\|_\infty := \max_{(i,j)} |Y_{ij}| = \max_{\|\mathbf{X}\|_1 \leq 1} \langle \mathbf{Y}, \mathbf{X} \rangle$ , the optimization problem

$$(2.20) \quad \text{minimize}_{\mathbf{X}} \quad \langle \mathbf{Y}, \mathbf{X} \rangle \quad \text{subject to } \|\mathbf{X}\|_1 \leq 1$$

has optimal value  $-\|\mathbf{Y}\|_\infty$ . One minimizer is the one-sparse matrix

$$\mathbf{X}^* = -\text{sgn}(Y_{i^*j^*})\mathbf{e}_{i^*}\mathbf{e}_{j^*}^\top,$$

where  $(i^*, j^*) \in \arg \max_{(i,j)} |Y_{ij}|$ ; i.e.  $\mathbf{X}^*$  has exactly one nonzero element.



**Projection onto the  $\ell_1$ -ball.** To effectively handle the sparse term in the norm constrained problem (1.6), we will need to modify the Frank-Wolfe algorithm by incorporating additional projection steps. For any  $\mathbf{Y} \in \mathbb{R}^{m \times n}$  and  $\beta > 0$ , the projection onto the  $\ell_1$ -ball:

$$(2.21) \quad \mathcal{P}_{\|\cdot\|_1 \leq \beta}[\mathbf{Y}] = \arg \min_{\|\mathbf{X}\|_1 \leq \beta} \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2,$$

can be easily solved with  $O(mn(\log m + \log n))$  cost [32]. Moreover, a divide and conquer algorithm, achieving linear cost in expectation to solve (2.21), has also been proposed in [32].

**Proximal mapping of  $\ell_1$  norm.** To effectively handle the sparse term arising in problem (1.4), we will need to modify the Frank-Wolfe algorithm by incorporating additional proximal steps. For any  $\mathbf{Y} \in \mathbb{R}^{m \times n}$  and  $\lambda > 0$ , the proximal mapping of  $\ell_1$  norm has the following closed-form expression

$$(2.22) \quad \mathcal{T}_\lambda[\mathbf{Y}] = \arg \min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{X}\|_1,$$

where  $\mathcal{T}_\lambda : \mathbb{R} \rightarrow \mathbb{R}$  denotes the soft-thresholding operator  $\mathcal{T}_\lambda(x) = \text{sgn}(x) \max\{|x| - \lambda, 0\}$ , and extension to matrices is obtained by applying the scalar operator  $\mathcal{T}_\lambda(\cdot)$  to each element.

## 4.4 范数约束问题的FW-P方法

- 考虑问题:

$$(3.1) \quad \min_{\mathbf{L}, \mathbf{S}} l(\mathbf{L}, \mathbf{S}) = \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 \quad \text{s.t.} \quad \|\mathbf{L}\|_* \leq \tau_L, \|\mathbf{S}\|_1 \leq \tau_S$$

注意到目标函数可微, 有:

$$(3.2) \quad \nabla_{\mathbf{L}} l(\mathbf{L}, \mathbf{S}) = \mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]$$

$$(3.3) \quad \nabla_{\mathbf{S}} l(\mathbf{L}, \mathbf{S}) = \mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]$$

$$\nabla l(\mathbf{L}, \mathbf{S}) = (\nabla_{\mathbf{L}} l, \nabla_{\mathbf{S}} l)$$

- 第k次FW算法产生的迭代点记为  $\mathbf{x}^k = (\mathbf{L}^k, \mathbf{S}^k)$ ,  $\mathbf{v}^k = (\mathbf{V}_L^k, \mathbf{V}_S^k)$  由下式得到:

$$(3.5) \quad \begin{pmatrix} \mathbf{V}_L^k \\ \mathbf{V}_S^k \end{pmatrix} \in \arg \min \left\langle \begin{pmatrix} \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \\ \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \end{pmatrix}, \begin{pmatrix} \mathbf{V}_L \\ \mathbf{V}_S \end{pmatrix} \right\rangle$$
$$\text{s.t.} \quad \|\mathbf{V}_L\|_* \leq \tau_L, \|\mathbf{V}_S\|_1 \leq \tau_S,$$

- (3.5) 是可分成两个独立的子问题:

$$\mathbf{V}_L^k \in \arg \min_{\|\mathbf{V}_L\|_* \leq \tau_L} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_L \rangle$$

$$\mathbf{V}_S^k \in \arg \min_{\|\mathbf{V}_S\|_1 \leq \tau_S} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_S \rangle.$$

- 容易解得:

$$(3.6) \quad \mathbf{V}_L^k = -\tau_L \mathbf{u}^k (\mathbf{v}^k)^\top,$$

$$(3.7) \quad \mathbf{V}_S^k = -\tau_S \cdot \delta_{i^* j^*}^k \cdot \mathbf{e}_{i^*}^k (\mathbf{e}_{j^*}^k)^\top$$

- 其中  $\mathbf{u}^k$  和  $\mathbf{v}^k$  是  $\mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]$  的对应最大奇异值的左右奇异向量。 $(i^*, j^*)$  是  $\mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]$  中最大元对应的位置。

$$\delta_{ij}^k := \text{sgn} \left[ (\mathcal{P}_Q [\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}])_{ij} \right]$$

---

**Algorithm 3** Frank-Wolfe method for problem (3.1)

---

```
1: Initialization:  $\mathbf{L}^0 = \mathbf{S}^0 = \mathbf{0}$ ;  
2: for  $k = 0, 1, 2, \dots$  do  
3:    $\mathbf{D}_L^k \in \arg \min_{\|\mathbf{D}_L\|_* \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L \rangle$ ;  $\mathbf{V}_L^k = \tau_L \mathbf{D}_L^k$ ;  
4:    $\mathbf{D}_S^k \in \arg \min_{\|\mathbf{D}_S\|_1 \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_S \rangle$ ;  $\mathbf{V}_S^k = \tau_S \mathbf{D}_S^k$ ;  
5:    $\gamma = \frac{2}{k+2}$ ;  
6:    $\mathbf{L}^{k+1} = \mathbf{L}^k + \gamma(\mathbf{V}_L^k - \mathbf{L}^k)$ ;  
7:    $\mathbf{S}^{k+1} = \mathbf{S}^k + \gamma(\mathbf{V}_S^k - \mathbf{S}^k)$ ;  
8: end for
```

---

- 不足:  $\mathbf{S}$  每次迭代只进行一次稀疏更新, 由于  $\mathbf{V}_S^k = -\tau_S \mathbf{e}_{i^*}^k (\mathbf{e}_{j^*}^k)^\top$  只有一个非零元, 实际应用中是一个明显的缺点, 作为最优的  $\mathbf{S}^*$  可能有相对较多的非零项

## 4.5 FW-P算法:将FW和投影梯度相结合

- 在第 $k$ 次迭代,  $(\mathbf{L}^{k+1/2}, \mathbf{S}^{k+1/2})$  由FW算法产生
- 对稀疏项  $\mathbf{S}$  增加一个额外的投影梯度步骤

$$(3.8) \quad \mathbf{S}^{k+1} = \mathcal{P}_{\|\cdot\|_1 \leq \tau_S} \left[ \mathbf{S}^{k+\frac{1}{2}} - \nabla_{\mathbf{S}} l(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}) \right]$$

$$(3.9) \quad = \mathcal{P}_{\|\cdot\|_1 \leq \tau_S} \left[ \mathbf{S}^{k+\frac{1}{2}} - \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}] \right]$$

---

**Algorithm 4** FW-P method for problem (3.1)

---

- 1: **Initialization:**  $\mathbf{L}^0 = \mathbf{S}^0 = \mathbf{0}$ ;
  - 2: **for**  $k = 0, 1, 2, \dots$  **do**
  - 3:    $\mathbf{D}_L^k \in \arg \min_{\|\mathbf{D}_L\|_* \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L \rangle$ ;  $\mathbf{V}_L^k = \tau_L \mathbf{D}_L^k$ ;
  - 4:    $\mathbf{D}_S^k \in \arg \min_{\|\mathbf{D}_S\|_1 \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_S \rangle$ ;  $\mathbf{V}_S^k = \tau_S \mathbf{D}_S^k$ ;
  - 5:    $\gamma = \frac{2}{k+2}$ ;
  - 6:    $\mathbf{L}^{k+\frac{1}{2}} = \mathbf{L}^k + \gamma(\mathbf{V}_L^k - \mathbf{L}^k)$ ;
  - 7:    $\mathbf{S}^{k+\frac{1}{2}} = \mathbf{S}^k + \gamma(\mathbf{V}_S^k - \mathbf{S}^k)$ ;
  - 8:    $\mathbf{S}^{k+1} = \mathcal{P}_{\|\cdot\|_1 \leq \tau_S} [\mathbf{S}^{k+\frac{1}{2}} - \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}]]$ ;
  - 9:    $\mathbf{L}^{k+1} = \mathbf{L}^{k+\frac{1}{2}}$ ;
  - 10: **end for**
-

## 4.6 使用合成数据的例子

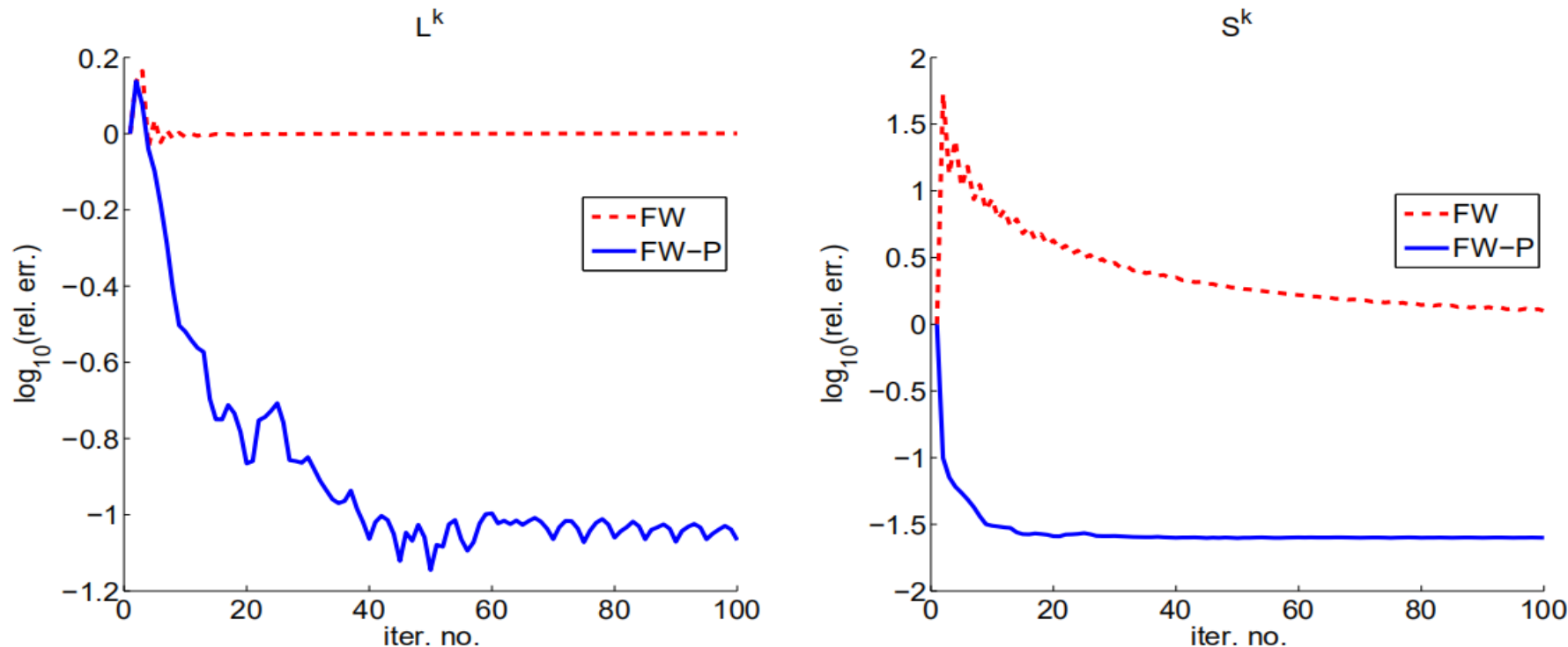


FIG. 1. *Comparisons between Algorithms 3 and 4 for problem (3.1) on synthetic data.* The data are generated in Matlab as  $m = 1000$ ;  $n = 1000$ ;  $r = 5$ ;  $L_0 = \text{randn}(m, r) * \text{randn}(r, n)$ ;  $\Omega_{\text{mask}} = \text{ones}(m, n)$ ;  $S_0 = 100 * \text{randn}(m, n) * (\text{rand}(m, n) < 0.01)$ ;  $M = L_0 + S_0 + \text{randn}(m, n)$ ;  $\tau_L = \text{norm\_nuc}(L_0)$ ;  $\tau_S = \text{norm}(\text{vec}(S_0), 1)$ ; The left figure plots  $\log_{10}(\|L^k - L_0\|_F / \|L_0\|_F)$  versus the iteration number  $k$ . The right figure plots  $\log_{10}(\|S^k - S_0\|_F / \|S_0\|_F)$  versus  $k$ . The FW-P method is clearly more efficient than the straightforward FW method in recovering  $L_0$  and  $S_0$ .

LEMMA 3.2. *The feasible set  $\mathcal{D} = \{(\mathbf{L}, \mathbf{S}) \mid \|\mathbf{L}\|_* \leq \tau_L, \|\mathbf{S}\|_1 \leq \tau_S\}$  has diameter  $D \leq 2\sqrt{\tau_L^2 + \tau_S^2}$ .*

*Proof.* For any  $\mathbf{Z} = (\mathbf{L}, \mathbf{S})$  and  $\mathbf{Z}' = (\mathbf{L}', \mathbf{S}') \in \mathcal{D}$ ,

$$\begin{aligned}
 \|\mathbf{Z} - \mathbf{Z}'\|_F^2 &= \|\mathbf{L} - \mathbf{L}'\|_F^2 + \|\mathbf{S} - \mathbf{S}'\|_F^2 \leq (\|\mathbf{L}\|_F + \|\mathbf{L}'\|_F)^2 + (\|\mathbf{S}\|_F + \|\mathbf{S}'\|_F)^2 \\
 (3.4) \quad &\leq (\|\mathbf{L}\|_* + \|\mathbf{L}'\|_*)^2 + (\|\mathbf{S}\|_1 + \|\mathbf{S}'\|_1)^2 \leq 4\tau_L^2 + 4\tau_S^2.
 \end{aligned}$$

$$\bullet \quad \|\mathbf{L}\|_F = \sqrt{\text{Tr}(\mathbf{L}^T \mathbf{L})} = \sqrt{\sum \sigma_i^2} \leq \sum \sigma_i = \|\mathbf{L}\|_*$$

## 4.7 FW-P算法的收敛性

**THEOREM 3.3.** *Let  $l^\star$  be the optimal value to problem (3.1),  $\mathbf{x}^k = (\mathbf{L}^k, \mathbf{S}^k)$  and  $\mathbf{v}^k = (\mathbf{V}_L^k, \mathbf{V}_S^k)$  be the sequence produced by Algorithm 4. Then we have*

$$(3.11) \quad l(\mathbf{L}^k, \mathbf{S}^k) - l^\star \leq \frac{16(\tau_L^2 + \tau_S^2)}{k + 2}.$$

*Moreover, for any  $K \geq 1$ , there exists  $1 \leq \tilde{k} \leq K$  such that the surrogate duality gap (defined in (2.13)) satisfies*

$$(3.12) \quad d(\mathbf{x}^{\tilde{k}}) = \left\langle \mathbf{x}^{\tilde{k}} - \mathbf{v}^{\tilde{k}}, \nabla l(\mathbf{x}^{\tilde{k}}) \right\rangle \leq \frac{48(\tau_L^2 + \tau_S^2)}{K + 2}.$$

*Proof.* Substituting  $L = 2$  (Lemma 3.1) and  $D \leq 2\sqrt{\tau_L^2 + \tau_S^2}$  (Lemma 3.2) into Theorems 2.1 and 2.2, we can easily obtain the above result.  $\square$



## 4.8 带有惩罚项的CPCP问题

- 考虑如下问题:

$$(4.1) \quad \min_{\mathbf{L}, \mathbf{S}} f(\mathbf{L}, \mathbf{S}) \doteq \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1$$

- 将上述问题转化为复合优化问题:

$$(4.2) \quad \begin{aligned} \min \quad & g(\mathbf{L}, \mathbf{S}, t_L, t_S) \doteq \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L t_L + \lambda_S t_S \\ \text{s.t.} \quad & \|\mathbf{L}\|_* \leq t_L, \quad \|\mathbf{S}\|_1 \leq t_S, \end{aligned}$$

- 注意到目标函数  $g(\mathbf{L}, \mathbf{S}, t_L, t_S)$  是可微的:

$$(4.3) \quad \nabla_{\mathbf{L}} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \nabla_{\mathbf{S}} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]$$

$$(4.4) \quad \nabla_{t_L} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \lambda_L, \quad \nabla_{t_S} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \lambda_S.$$

梯度  $\nabla g(\mathbf{L}, \mathbf{S}, t_L, t_S) = (\nabla_{\mathbf{L}} g, \nabla_{\mathbf{S}} g, \nabla_{t_L} g, \nabla_{t_S} g)$

- 但是此时Frank-Wolfe方法仍然不能处理，因为可行域是无界。如果能得到 $t_L$ 和 $t_S$ 的最优值的上界即可利用FW方法。此时有：

$$(4.6) \quad \begin{aligned} \min \quad & \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L t_L + \lambda_S t_S \\ \text{s.t.} \quad & \|\mathbf{L}\|_* \leq t_L \leq U_L, \quad \|\mathbf{S}\|_1 \leq t_S \leq U_S, \end{aligned}$$

- 注意到 $\mathbf{L} = \mathbf{0}, \mathbf{S} = \mathbf{0}, t_L = 0, t_S = 0$  时，目标函数值为 $\frac{1}{2} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2$ ，可令：

$$(4.7) \quad t_L^* \leq \frac{1}{2\lambda_L} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2, \quad t_S^* \leq \frac{1}{2\lambda_S} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2$$

LEMMA 4.2. *The feasible set  $\mathcal{D} = \{(\mathbf{L}, \mathbf{S}, t_L, t_S) \mid \|\mathbf{L}\|_* \leq t_L \leq U_L, \|\mathbf{S}\|_1 \leq t_S \leq U_S\}$  has diameter  $D \leq \sqrt{5} \cdot \sqrt{U_L^2 + U_S^2}$ .*

*Proof.* Since for any  $\mathbf{Z} = (\mathbf{L}, \mathbf{S}, t_L, t_S), \mathbf{Z}' = (\mathbf{L}', \mathbf{S}', t'_L, t'_S) \in \mathcal{D}$ , we have

$$\begin{aligned} \|\mathbf{Z} - \mathbf{Z}'\|_F^2 &= \|\mathbf{L} - \mathbf{L}'\|_F^2 + \|\mathbf{S} - \mathbf{S}'\|_F^2 + (t_L - t'_L)^2 + (t_S - t'_S)^2 \\ &\leq (\|\mathbf{L}\|_F + \|\mathbf{L}'\|_F)^2 + (\|\mathbf{S}\|_F + \|\mathbf{S}'\|_F)^2 + (t_L - t'_L)^2 + (t_S - t'_S)^2 \\ &\leq (\|\mathbf{L}\|_* + \|\mathbf{L}'\|_*)^2 + (\|\mathbf{S}\|_1 + \|\mathbf{S}'\|_1)^2 + (t_L - t'_L)^2 + (t_S - t'_S)^2 \\ &\leq (U_L + U_L)^2 + (U_S + U_S)^2 + U_L^2 + U_S^2 \\ &= 5(U_L^2 + U_S^2) \end{aligned}$$

## 4.9 问题(4.6)的FW算法

- 第k步, 迭代点  $\mathbf{x}^k = (\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)$  由FW算法产生  
令  $\mathbf{v}^k = (\mathbf{V}_L^k, \mathbf{V}_S^k, V_{t_L}^k, V_{t_S}^k)$  是求解下述问题(FW)得到的:

$$(4.9) \quad \mathbf{v}^k \in \arg \min_{\mathbf{v} \in \mathcal{D}} \quad \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_L + \mathbf{V}_S \rangle + \lambda_L V_{t_L} + \lambda_S V_{t_S}$$

- 上述问题可以解耦成两个独立的子问题:

$$(4.10) \quad (\mathbf{V}_L^k, V_{t_L}^k) \in \arg \min_{\|\mathbf{V}_L\|_* \leq V_{t_L} \leq U_L} g_L(\mathbf{V}_L, V_{t_L}) \doteq \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_L \rangle + \lambda_L V_{t_L}$$

$$(4.11) \quad (\mathbf{V}_S^k, V_{t_S}^k) \in \arg \min_{\|\mathbf{V}_S\|_1 \leq V_{t_S} \leq U_S} g_S(\mathbf{V}_S, V_{t_S}) \doteq \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_S \rangle + \lambda_S V_{t_S}.$$

- 先考虑第一个子问题, 令:

$$(4.12) \quad \mathbf{D}_L^k \in \arg \min_{\|\mathbf{D}_L\|_* \leq 1} \hat{g}_L(\mathbf{D}_L) \doteq \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L \rangle + \lambda_L$$

- 利用[24,引理1]:

$$(4.13) \quad (\mathbf{V}_L^k, V_{t_L}^k) \in \begin{cases} \{(\mathbf{0}, 0)\} & \text{if } \hat{g}_L(\mathbf{D}_L^k) > 0 \\ \text{conv}\{(\mathbf{0}, 0), U_L(\mathbf{D}_L^k, 1)\} & \text{if } \hat{g}_L(\mathbf{D}_L^k) = 0 \\ \{U_L(\mathbf{D}_L^k, 1)\} & \text{if } \hat{g}_L(\mathbf{D}_L^k) < 0. \end{cases}$$

COROLLARY 4.3. *Let  $\mathbf{x}^* = (\mathbf{L}^*, \mathbf{S}^*, t_L^*, t_S^*)$  be an optimal solution to (4.6). For  $\{\mathbf{x}^k\}$  generated by Algorithm 5, we have for  $k = 0, 1, 2, \dots$ ,*

$$(4.14) \quad g(\mathbf{x}^k) - g(\mathbf{x}^*) \leq \frac{20(U_L^2 + U_S^2)}{k + 2}.$$

[24] Z. Harchaoui, A. Juditsky, and A. Nemirovski, “Conditional gradient algorithms for norm-regularized smooth convex optimization,” *Mathematical Programming*, pp. 1–38, 2014

*Proof.* Applying Theorem 2.1 with parameters calculated in Lemmas 4.1 and 4.2 we directly have

$$(4.15) \quad g(\mathbf{x}^k) - g(\mathbf{x}^\star) \leq \frac{2 \cdot 2 \cdot \left( \sqrt{5(U_L^2 + U_S^2)} \right)^2}{k + 2} = \frac{20(U_L^2 + U_S^2)}{k + 2}.$$

A more careful calculation below slightly improves the constant in (4.15).

$$(4.16) \quad \begin{aligned} g(\mathbf{x}^{k+1}) &= g(\mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k)) \\ &\leq g(\mathbf{x}^k) + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + \gamma^2 \|\mathbf{V}_L^k - \mathbf{L}^k\|_F^2 + \gamma^2 \|\mathbf{V}_S^k - \mathbf{S}^k\|_F^2 \\ &\leq g(\mathbf{x}^k) + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + 4\gamma^2(U_L^2 + U_S^2), \end{aligned}$$

where the second line holds by noting that  $g$  is only linear in  $t_L$  and  $t_S$ ; the last line holds as

$$\begin{aligned} \|\mathbf{V}_L^k - \mathbf{L}^k\|_F^2 &\leq (\|\mathbf{V}_L^k\|_F + \|\mathbf{L}^k\|_F)^2 \leq (U_L + U_L)^2 = 4U_L^2, \quad \text{and} \\ \|\mathbf{V}_S^k - \mathbf{S}^k\|_F^2 &\leq (\|\mathbf{V}_S^k\|_F + \|\mathbf{S}^k\|_F)^2 \leq (U_S + U_S)^2 = 4U_S^2. \end{aligned}$$

Following the arguments in the proof of Theorem 1 with (2.10) replaced by (4.16), we can easily obtain that

$$g(\mathbf{x}^k) - g(\mathbf{x}^\star) \leq \frac{16(U_L^2 + U_S^2)}{k + 2}.$$

---

**Algorithm 5** Frank-Wolfe method for problem (4.6)

---

```
1: Initialization:  $\mathbf{L}^0 = \mathbf{S}^0 = \mathbf{0}$ ;  $t_L^0 = t_S^0 = 0$ ;  
2: for  $k = 0, 1, 2, \dots$  do  
3:    $\mathbf{D}_L^k \in \arg \min_{\|\mathbf{D}_L\|_* \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L \rangle$ ;  
4:    $\mathbf{D}_S^k \in \arg \min_{\|\mathbf{D}_S\|_1 \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_S \rangle$ ;  
5:   if  $\lambda_L \geq -\langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L^k \rangle$  then  
6:      $\mathbf{V}_L^k = \mathbf{0}$ ;  $V_{t_L}^k = 0$   
7:   else  
8:      $\mathbf{V}_L^k = U_L \mathbf{D}_L^k$ ,  $V_{t_L}^k = U_L$ ;  
9:   end if  
10:  if  $\lambda_S \geq -\langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_S^k \rangle$  then  
11:     $\mathbf{V}_S^k = \mathbf{0}$ ;  $V_{t_S}^k = 0$ ;  
12:  else  
13:     $\mathbf{V}_S^k = U_S \mathbf{D}_S^k$ ,  $V_{t_S}^k = U_S$ ;  
14:  end if  
15:   $\gamma = \frac{2}{k+2}$ ;  
16:   $\mathbf{L}^{k+1} = (1 - \gamma)\mathbf{L}^k + \gamma\mathbf{V}_L^k$ ,  $t_L^{k+1} = (1 - \gamma)t_L^k + \gamma V_{t_L}^k$ ;  
17:   $\mathbf{S}^{k+1} = (1 - \gamma)\mathbf{S}^k + \gamma\mathbf{V}_S^k$ ,  $t_S^{k+1} = (1 - \gamma)t_S^k + \gamma V_{t_S}^k$ ;  
18: end for
```

---

**算法5的不足：**首先在范数约束的情况下， $\mathbf{S}$  只有一个非零项，所以  $\mathbf{V}_S^k$  在每次迭代中只有一个稀疏更新。其次，精确的收敛速度依赖于我们对  $U_L$  和  $U_S$  的估计。



## 4.10 问题(4.6)的FW-T算法:结合FW和临近梯度法

- 主要思想:利用FW算法更新 $\mathbf{L}$ , 利用临近梯度法更新  $\mathbf{S}$ , 精确线搜索, 自适应更新 $U_L$ 和 $U_S$

*Proximal gradient step for  $\mathbf{S}$ .* To update  $\mathbf{S}$  in a more efficient way, we incorporate an additional proximal gradient step for  $\mathbf{S}$ . At iteration  $k$ , let  $(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}})$  be the result produced by Frank-Wolfe step. To produce the next iterate, we retain the low-rank term  $\mathbf{L}^{k+\frac{1}{2}}$ , but execute a proximal gradient step for the function  $f(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S})$  at the point  $\mathbf{S}^{k+\frac{1}{2}}$ , i.e.

$$\begin{aligned} \mathbf{S}^{k+1} &\in \arg \min_{\mathbf{S}} \left\langle \nabla_{\mathbf{S}} f(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}), \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\rangle + \frac{1}{2} \left\| \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\|_F^2 + \lambda_S \|\mathbf{S}\|_1 \\ (4.17) \quad &= \arg \min_{\mathbf{S}} \left\langle \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}], \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\rangle + \frac{1}{2} \left\| \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\|_F^2 + \lambda_S \|\mathbf{S}\|_1 \end{aligned}$$

which can be easily computed using the soft-thresholding operator:

$$(4.18) \quad \mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S} \left[ \mathbf{S}^{k+\frac{1}{2}} - \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}] \right].$$

**Exact line search.** For the Frank-Wolfe step, instead of choosing the fixed step length  $\frac{2}{k+2}$ , we implement an exact line search by solving a two-dimensional quadratic problem (4.20), as in [24]. This modification turns out to be crucial to achieve a primal convergence result that only weakly depends on the tightness of our guesses  $U_L$  and  $U_S$ .

**Adaptive updates of  $U_L$  and  $U_S$ .** We initialize  $U_L$  and  $U_S$  using the crude bound (4.8). Then, at the end of the  $k$ -iteration, we respectively update

$$(4.19) \quad U_L^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_L, \quad U_S^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_S.$$

This scheme maintains the property that  $U_L^{k+1} \geq t_L^*$  and  $U_S^{k+1} \geq t_S^*$ . Moreover, we prove (Lemma 4.4) that  $g$  is non-increasing through our algorithm, and so this scheme produces a sequence of tighter upper bounds for  $U_L^*$  and  $U_S^*$ . Although this dynamic scheme does not improve the theoretical convergence result, some acceleration is empirically exhibited.



---

**Algorithm 6** FW-T method for problem (4.1)

---

- 1: **Input:** data matrix  $\mathbf{M} \in \mathbb{R}^{m \times n}$ ; weights  $\lambda_L, \lambda_S > 0$ ; max iteration number  $T$ ;
- 2: **Initialization:**  $\mathbf{L}^0 = \mathbf{S}^0 = \mathbf{0}$ ;  $t_L^0 = t_S^0 = 0$ ;  $U_L^0 = g(\mathbf{L}^0, \mathbf{S}^0, t_L^0, t_S^0)/\lambda_L$ ;  $U_S^0 = g(\mathbf{L}^0, \mathbf{S}^0, t_L^0, t_S^0)/\lambda_S$ ;
- 3: **for**  $k = 0, 1, 2, \dots, T$  **do**
- 4:   *same as lines 3-14 in Algorithm 5;*
- 5:    $\left(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, t_S^{k+\frac{1}{2}}\right)$  is computed as an optimizer to

$$(4.20) \quad \begin{aligned} \min \quad & \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L t_L + \lambda_S t_S \\ \text{s.t.} \quad & \begin{pmatrix} \mathbf{L} \\ t_L \end{pmatrix} \in \text{conv} \left\{ \begin{pmatrix} \mathbf{L}^k \\ t_L^k \end{pmatrix}, \begin{pmatrix} \mathbf{V}_L^k \\ V_{t_L}^k \end{pmatrix} \right\} \\ & \begin{pmatrix} \mathbf{S} \\ t_S \end{pmatrix} \in \text{conv} \left\{ \begin{pmatrix} \mathbf{S}^k \\ t_S^k \end{pmatrix}, \begin{pmatrix} \mathbf{V}_S^k \\ V_{t_S}^k \end{pmatrix} \right\}; \end{aligned}$$

- 6:    $\mathbf{S}^{k+1} = \mathcal{T}[\mathbf{S}^{k+\frac{1}{2}} - \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}], \lambda_S]$ ;
  - 7:    $\mathbf{L}^{k+1} = \mathbf{L}^{k+\frac{1}{2}}, t_L^{k+1} = t_L^{k+\frac{1}{2}}; t_S^{k+1} = \|\mathbf{S}^{k+1}\|_1$ ;
  - 8:    $U_L^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_L$ ;
  - 9:    $U_S^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_S$ ;
  - 10: **end for**
-

## 4.11 FW-T算法的收敛性

THEOREM 4.5. *Let  $r_L^*$  and  $r_S^*$  be the smallest radii such that*

$$(4.24) \quad \left\{ (\mathbf{L}, \mathbf{S}) \mid f(\mathbf{L}, \mathbf{S}) \leq g(\mathbf{L}^0, \mathbf{S}^0, t_L^0, t_S^0) = \frac{1}{2} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2 \right\} \subseteq \overline{B(r_L^*)} \times \overline{B(r_S^*)},$$

where  $\overline{B(r)} \doteq \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \|\mathbf{X}\|_F \leq r\}$  for any  $r \geq 0$ .<sup>‡</sup> Then for the sequence  $\{(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)\}$  generated by Algorithm 6, we have

$$(4.25) \quad g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k) - g(\mathbf{L}^*, \mathbf{S}^*, t_L^*, t_S^*) \\ \leq \frac{\min\{4(t_L^* + r_L^*)^2 + 4(t_S^* + r_S^*)^2, 16(U_L^0)^2 + 16(U_S^0)^2\}}{k + 2}.$$

THEOREM 4.6. Let  $\mathbf{x}^k$  denote  $(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)$  generated by Algorithm 6. Then for any  $K \geq 1$ , there exists  $1 \leq \tilde{k} \leq K$  such that

$$(4.29) \quad g(\mathbf{x}^{\tilde{k}}) - g(\mathbf{x}^*) \leq d(\mathbf{x}^{\tilde{k}}) \leq \frac{48 ((U_L^0)^2 + (U_S^0)^2)}{K + 2}.$$

*Proof.* Define  $\Delta^k = g(\mathbf{x}^k) - g(\mathbf{x}^*)$ . Following (4.16), we have

$$(4.30) \quad \Delta^{k+1} \leq \Delta^k + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + 4\gamma^2 ((U_L^0)^2 + (U_S^0)^2).$$

Then following the arguments in the proof of Theorem 2 with (2.17) replaced by (4.30), we can easily obtain the result.  $\square$

**Stopping criterion.** Compared to the convergence of  $g(\mathbf{x}^k)$  (Theorem 4.5), the convergence result for  $d(\mathbf{x}^k)$  can be much slower (Theorem 4.6). Therefore, here the surrogate duality gap  $d(\cdot)$  is not that suitable to serve as a stopping criterion. Consequently, in our implementation, we terminate Algorithm 6 if

$$(4.31) \quad |g(\mathbf{x}^{k+1}) - g(\mathbf{x}^k)| / g(\mathbf{x}^k) \leq \varepsilon,$$

for five consecutive iterations.

## 5. 数值实验

### 5.1 ISTA算法与FISTA算法

- ISTA算法在第k步通过求解下式来更新  $(\mathbf{L}, \mathbf{S})$  :

$$\begin{aligned} (\mathbf{L}^{k+1}, \mathbf{S}^{k+1}) = \arg \min_{\mathbf{L}, \mathbf{S}} & \left\langle \begin{pmatrix} \nabla_{\mathbf{L}} l(\mathbf{L}^k, \mathbf{S}^k) \\ \nabla_{\mathbf{S}} l(\mathbf{L}^k, \mathbf{S}^k) \end{pmatrix}, \begin{pmatrix} \mathbf{L} - \mathbf{L}^k \\ \mathbf{S} - \mathbf{S}^k \end{pmatrix} \right\rangle + \\ (5.1) \quad & \frac{L_f}{2} \left\| \begin{pmatrix} \mathbf{L} \\ \mathbf{S} \end{pmatrix} - \begin{pmatrix} \mathbf{L}^k \\ \mathbf{S}^k \end{pmatrix} \right\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1 \end{aligned}$$

Here  $L_f = 2$  denotes the Lipschitz constant of  $\nabla l(\mathbf{L}, \mathbf{S})$  with respect to  $(\mathbf{L}, \mathbf{S})$ , and  $\nabla_{\mathbf{L}} l(\mathbf{L}^k, \mathbf{S}^k) = \nabla_{\mathbf{S}} l(\mathbf{L}^k, \mathbf{S}^k) = \mathcal{P}_{\Omega}[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]$ . Since  $\mathbf{L}$  and  $\mathbf{S}$  are decoupled in (5.1), equivalently we have

$$(5.2) \quad \mathbf{L}^{k+1} = \arg \min_{\mathbf{L}} \left\| \mathbf{L} - \left( \mathbf{L}^k - \frac{1}{2} \mathcal{P}_{\Omega}[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right) \right\|_F^2 + \lambda_L \|\mathbf{L}\|_*,$$

$$(5.3) \quad \mathbf{S}^{k+1} = \arg \min_{\mathbf{S}} \left\| \mathbf{S} - \left( \mathbf{S}^k - \frac{1}{2} \mathcal{P}_{\Omega}[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right) \right\|_F^2 + \lambda_S \|\mathbf{S}\|_1.$$

The solution to problem (5.3) can be given explicitly in terms of the proximal mapping of  $\|\cdot\|_1$  as introduced in Section 2.2, i.e.,

$$\mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S/2} \left[ \mathbf{S}^k - \frac{1}{2} \mathcal{P}_{\Omega}[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right].$$

For a matrix  $\mathbf{X}$  and any  $\tau \geq 0$ , let  $\mathcal{D}_{\tau}(\mathbf{X})$  denote the singular value thresholding operator  $\mathcal{D}_{\tau}(\mathbf{X}) = \mathbf{U} \mathcal{T}_{\tau}(\mathbf{\Sigma}) \mathbf{V}^{\top}$ , where  $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top}$  is the singular value decomposition of  $\mathbf{X}$ . It is not difficult to show [35, 36] that the solution to problem (5.2) can be given explicitly by

$$\mathbf{L}^{k+1} = \mathcal{D}_{\lambda_L/2} \left[ \mathbf{L}^k - \frac{1}{2} \mathcal{P}_{\Omega}[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right].$$

---

**Algorithm 7** ISTA for problem (1.5)

---

1: **Initialization:**  $\mathbf{L}^0 = \mathbf{0}$ ,  $\mathbf{S}^0 = \mathbf{0}$ ;  
2: **for**  $k = 0, 1, 2, \dots$  **do**  
3:    $\mathbf{L}^{k+1} = \mathcal{D}_{\lambda_L/2} \left[ \mathbf{L}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right]$ ;  
4:    $\mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S/2} \left[ \mathbf{S}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right]$ ;  
5: **end for**

---

---

**Algorithm 8** FISTA for problem (1.5)

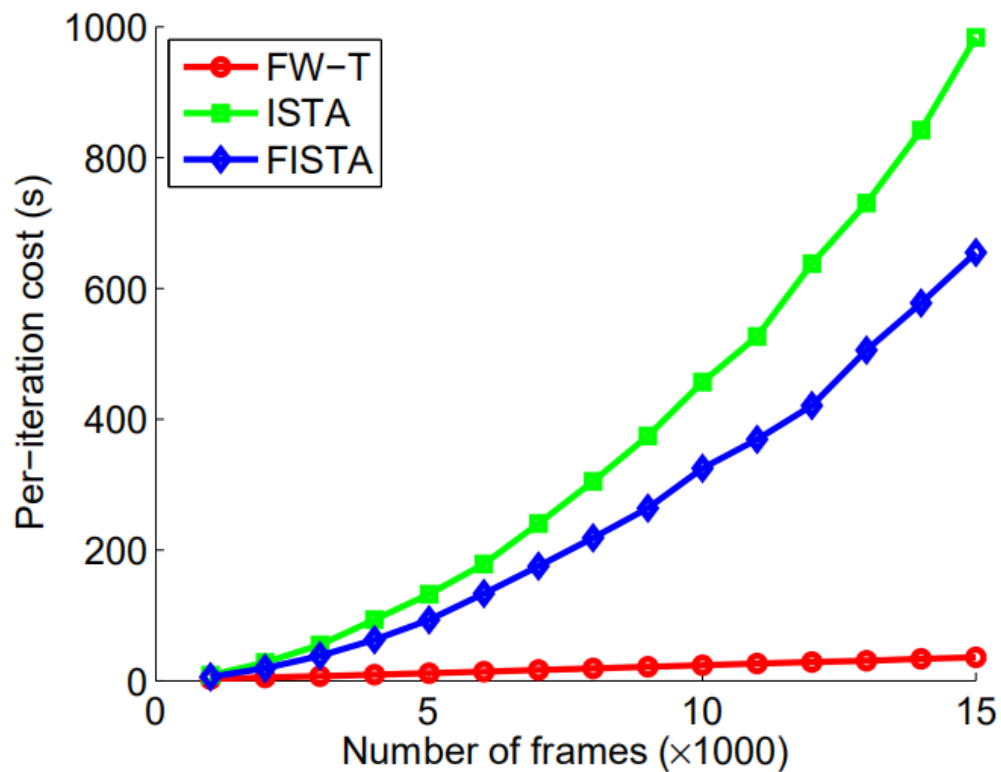
---

1: **Initialization:**  $\hat{\mathbf{L}}^0 = \mathbf{L}^0 = \mathbf{0}$ ,  $\hat{\mathbf{S}}^0 = \mathbf{S}^0 = \mathbf{0}$ ,  $t_0 = 1$ ;  
2: **for**  $k = 0, 1, 2, \dots$  **do**  
3:    $\mathbf{L}^{k+1} = \mathcal{D}_{\lambda_L/2} \left[ \hat{\mathbf{L}}^k - \frac{1}{2} \mathcal{P}_\Omega[\hat{\mathbf{L}}^k + \hat{\mathbf{S}}^k - \mathbf{M}] \right]$ ;  
4:    $\mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S/2} \left[ \hat{\mathbf{S}}^k - \frac{1}{2} \mathcal{P}_\Omega[\hat{\mathbf{L}}^k + \hat{\mathbf{S}}^k - \mathbf{M}] \right]$ ;  
5:    $t^{k+1} = \frac{1 + \sqrt{1 + 4(t^k)^2}}{2}$ ;  
6:    $\hat{\mathbf{L}}^{k+1} = \mathbf{L}^{k+1} + \frac{t^k - 1}{t^{k+1}} (\mathbf{L}^{k+1} - \mathbf{L}^k)$ ;  
7:    $\hat{\mathbf{S}}^{k+1} = \mathbf{S}^{k+1} + \frac{t^k - 1}{t^{k+1}} (\mathbf{S}^{k+1} - \mathbf{S}^k)$ ;  
8: **end for**

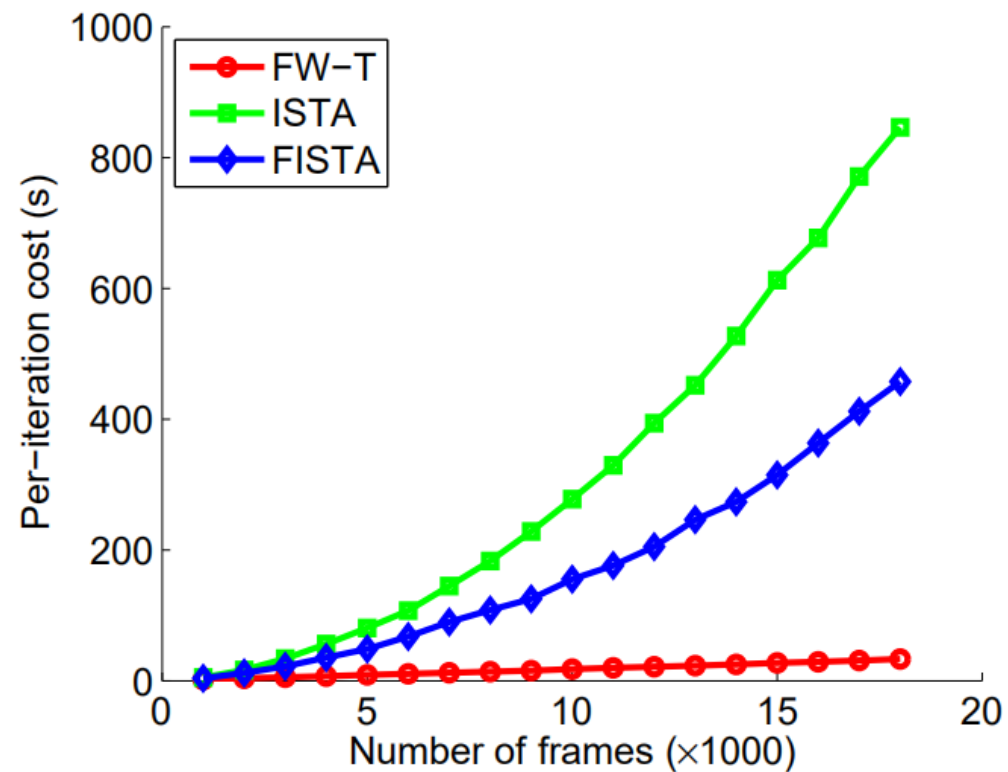
---



## 5.2 ISTA算法与FISTA算法



Airport



Square

FIG. 2. *Per-iteration cost vs. the number of frames in Airport and Square videos with full observation.* The per-iteration cost of our FW-T method grows linearly with the size of data, in contrast with the superlinear per-iteration cost of ISTA and FISTA. That makes the FW-T method more advantageous or may even be the only feasible choice for large problems.

TABLE 1

*Comparisons of FW-T, ISTA and FISTA on surveillance video data. The advantage of our FW-T method becomes prominent when the data are at large scale and compressed (i.e. the small  $\rho$  scenarios).*

Data	$\rho$	FW-T		ISTA		FISTA	
		iter.	time	iter.	time	iter.	time
Lobby (20480 $\times$ 1000)	1.0	96	1.94e+02	144	3.64e+02	41	<b>1.60e+02</b>
	0.8	104	<b>2.33e+02</b>	216	1.03e+03	52	3.55e+02
	0.6	133	<b>3.12e+02</b>	380	1.67e+03	74	5.10e+02
Campus (20480 $\times$ 1439)	1.0	45	<b>1.56e+02</b>	78	1.49e+03	23	4.63e+02
	0.8	44	<b>1.57e+02</b>	122	2.34e+03	30	6.45e+02
	0.6	41	<b>1.39e+02</b>	218	4.27e+03	43	1.08e+03
Escalator (20800 $\times$ 3417)	1.0	81	<b>7.40e+02</b>	58	4.19e+03	25	2.18e+03
	0.8	80	<b>7.35e+02</b>	90	8.18e+03	32	3.46e+03
	0.6	82	<b>7.68e+02</b>	162	1.83e+04	43	5.73e+03
Mall (81920 $\times$ 1286)	1.0	38	<b>4.70e+02</b>	110	5.03e+03	35	1.73e+03
	0.8	35	<b>4.58e+02</b>	171	7.32e+03	44	2.34e+03
	0.6	44	<b>5.09e+02</b>	308	1.31e+04	62	3.42e+03
Restaurant (19200 $\times$ 3055)	1.0	70	<b>5.44e+02</b>	52	3.01e+03	20	1.63e+03
	0.8	74	<b>5.51e+02</b>	81	4.84e+03	26	1.82e+03
	0.6	76	<b>5.73e+02</b>	144	9.93e+03	38	3.31e+03



Hall (25344 × 3584)	1.0	60	<b>6.33e+02</b>	52	2.98e+03	21	1.39e+03
	0.8	62	<b>6.52e+02</b>	81	6.45e+03	28	2.90e+03
	0.6	70	<b>7.43e+02</b>	144	1.42e+04	39	4.94e+03
Airport (25344 × 15730)	1.0	130	<b>6.42e+03</b>	29	2.37e+04	14	1.37e+04
	0.8	136	<b>6.65e+03</b>	45	6.92e+04	18	4.27e+04
	0.6	154	<b>7.72e+03</b>	77	1.78e+05	24	7.32e+04
Square (19200 × 28181)	1.0	179	<b>1.24e+04</b>	29	3.15e+04	13	1.51e+04
	0.8	181	<b>1.26e+04</b>	44	1.04e+05	17	6.03e+04
	0.6	191	<b>1.31e+04</b>	78	2.63e+05	22	9.88e+05

**谢谢！**