

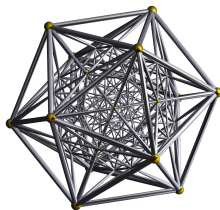
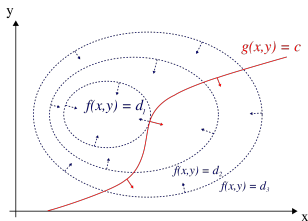
# Computational Principles for High-dim Data Analysis

## (Lecture Fourteen)

Yi Ma

EECS Department, UC Berkeley

October 14, 2021



# Constrained Convex Optimization for Structured Data Recovery

- 1 Constrained Optimization
- 2 Augmented Lagrangian Multipliers
- 3 Alternating Direction Method of Multipliers
- 4 More Scalable Algorithms

*"Since the fabric of the universe is most perfect and the work of a most wise Creator, nothing at all takes place in the universe in which some rule of maximum or minimum does not appear."*

– Leonhard Euler

# Optimization Challenges for Structured Data Recovery

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \doteq \underbrace{f(\mathbf{x})}_{\text{smooth convex}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth convex}}. \quad (1)$$

- **Challenge of Scale:** scale algorithms to when  $n$  is very large.

$$\text{Second order methods} \implies \text{First order methods...} \quad (2)$$

- **Nonsmoothness:** first order methods are slow for nonsmooth.

$$O(1/\sqrt{k}) \implies O(1/k) \implies O(1/k^2)... \quad (3)$$

- **Equality Constraints:** augmented Lagrange multiplier (ALM).
- **Separable Structures:** alternating direction of multipliers method (ADMM).

# Linear Equality Constrained Optimization

## Problem:

$$\min_{\mathbf{x}} g(\mathbf{x}) \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{y}, \quad (4)$$

where

- $g : \mathbb{R}^n \rightarrow \mathbb{R}$  is a (probably nonsmooth) convex function,
- $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{y} \in \text{range}(\mathbf{A})$  (so that the problem is feasible).

**A Natural Attempt:** solve the unconstrained by penalizing the constraint:

$$\hat{\mathbf{x}}(\mu) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 \quad \text{for a large } \mu. \quad (5)$$

- **Pros:** As  $\mu \rightarrow +\infty$ ,  $\hat{\mathbf{x}}(\mu) \rightarrow \mathbf{x}_\star$  (the “continuation method”).
- **Cons:** The rate of convergence depends on  $L = \mu \|\mathbf{A}\|_2^2$ .

# Lagrange Multiplier Method

## A More Principled Approach:

### Definition (The Lagrange Duality)

The *Lagrangian* function of the constrained problem (4):

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \doteq g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle, \quad (6)$$

where  $\boldsymbol{\lambda} \in \mathbb{R}^m$  is a vector of *Lagrange multipliers*. This gives a *dual function*:

$$d(\boldsymbol{\lambda}) \doteq \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle. \quad (7)$$

**Fact** (credited to Lagrange):  $\exists \boldsymbol{\lambda}_\star$  such that the optimal solution  $(\mathbf{x}_\star, \boldsymbol{\lambda}_\star)$  is a saddle point of the Lagrangian:

$$\sup_{\boldsymbol{\lambda}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \sup_{\boldsymbol{\lambda}} \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle = \sup_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda}). \quad (8)$$

# Dual Ascent Algorithm for the Lagrangian

**Fact:** If

$$\mathbf{x}'(\boldsymbol{\lambda}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle,$$

then  $\mathbf{A}\mathbf{x}'(\boldsymbol{\lambda}) - \mathbf{y}$  is a supergradient  $\partial d(\boldsymbol{\lambda})$  of the concave dual  $d(\boldsymbol{\lambda})$  at  $\boldsymbol{\lambda}$ .

(**Why?** Actually this is true for the dual function of general constraints  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ :  $d(\boldsymbol{\lambda}) = \min_{\mathbf{x}} g(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x})$ .)

**A Natural Attempt** to find the saddle point  $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$  is via *dual ascent*:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (9)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + t_{k+1}(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}). \quad (10)$$

- For certain problem classes, this converges to the optimal  $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$ .
- However, unfortunately it fails for problems in our settings.

# An Example of Failure

Consider the basis pursuit problem:

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{y}. \quad (11)$$

We can show that

$$\inf_{\mathbf{x}} \|\mathbf{x}\|_1 + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle = \begin{cases} -\infty & \|\mathbf{A}^* \boldsymbol{\lambda}\|_{\infty} > 1, \\ -\langle \boldsymbol{\lambda}, \mathbf{y} \rangle & \|\mathbf{A}^* \boldsymbol{\lambda}\|_{\infty} \leq 1. \end{cases} \quad (12)$$

Whenever the dual ascent step (10) happens to produce a  $\boldsymbol{\lambda}$  such that  $\|\mathbf{A}^* \boldsymbol{\lambda}\|_{\infty} > 1$ , the algorithm will break down.

**The reason is  $g(\mathbf{x})$  is not “strongly” convex enough.**

# Augmented Lagrange Multiplier

**One way out:** combining (5) and (4), consider the Augmented Lagrangian [Hestenes'69, Powell'69]:

$$\mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}) \doteq g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2. \quad (13)$$

Can be regarded as the Lagrangian for the penalized constrained problem

$$\min_{\mathbf{x}} g(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{y}, \quad (14)$$

which has the same optimal solution as the un-penalized problem.



# Augmented Lagrange Multiplier

Apply dual ascent to  $\mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda})$  with a particular step size  $t_{k+1} = \mu$ ,

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (15)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}). \quad (16)$$

**Fact:**  $\mathbf{x}_{k+1}$  always minimizes the unaugmented Lagrangian  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_{k+1})$  at  $\boldsymbol{\lambda} = \boldsymbol{\lambda}_{k+1}$ , because:

$$\begin{aligned} \mathbf{0} &\in \partial \mathcal{L}_\mu(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_k), \\ &= \partial g(\mathbf{x}_{k+1}) + \mathbf{A}^* \boldsymbol{\lambda}_k + \mu \mathbf{A}^* (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}), \\ &= \partial g(\mathbf{x}_{k+1}) + \mathbf{A}^* \boldsymbol{\lambda}_{k+1}, \\ &= \partial \mathcal{L}(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1}). \end{aligned}$$

**$\boldsymbol{\lambda}_{k+1}$  is always feasible, no bad behaviors!**

# Augmented Lagrange Multiplier

## Augmented Lagrange Multiplier (ALM)

**Problem Class:**  $\min_{\mathbf{x}} g(\mathbf{x})$  subject to  $\mathbf{Ax} = \mathbf{y}$ .  
 with  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  convex and coercive,  $\mathbf{y} \in \text{range}(\mathbf{A})$ .

**Basic Iteration:** set

$$\mathcal{L}_{\mu}(\mathbf{x}, \boldsymbol{\lambda}) = g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2.$$

Repeat:

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_{\mu}(\mathbf{x}, \boldsymbol{\lambda}_k),$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{Ax}_{k+1} - \mathbf{y}).$$

**Convergence Guarantee:**

$\{\mathbf{x}_k\}$  converges to an optimal solution at a rate  $O(1/k)$ .

# ALM for Basis Pursuit

## Augmented Lagrange Multiplier (ALM) for BP

- 1: **Problem:**  $\min_{\mathbf{x}} \|\mathbf{x}\|_1$  subject to  $\mathbf{y} = \mathbf{A}\mathbf{x}$ ,  
given  $\mathbf{y} \in \mathbb{R}^m$  and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\boldsymbol{\lambda}_0 \in \mathbb{R}^m$ , and  $\beta > 1$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:    $\mathbf{x}_{k+1} \leftarrow \arg \min \mathcal{L}_{\mu_k}(\mathbf{x}, \boldsymbol{\lambda}_k)$  using APG.
- 5:    $\boldsymbol{\lambda}_{k+1} \leftarrow \boldsymbol{\lambda}_k + \mu_k(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y})$ .
- 6:    $\mu_{k+1} \leftarrow \min\{\beta\mu_k, \mu_{\max}\}$ .
- 7: **end for**
- 8: **Output:**  $\mathbf{x}_\star \leftarrow \mathbf{x}_K$ .

# ALM for Principal Component Pursuit

## Augmented Lagrange Multiplier (ALM) for PCP

- 1: **Problem:**  $\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1$  subject to  $\mathbf{L} + \mathbf{S} = \mathbf{Y}$ ,  
given  $\mathbf{Y}$  and  $\lambda > 0$ .
- 2: **Input:**  $\mathbf{L}_0, \mathbf{S}_0, \mathbf{\Lambda}_0 \in \mathbb{R}^{m \times n}$  and  $\beta > 1$ .
- 3: **for**  $(k = 0, 1, 2, \dots, K - 1)$  **do**
- 4:    $\{\mathbf{L}_{k+1}, \mathbf{S}_{k+1}\} \leftarrow \arg \min \mathcal{L}_{\mu_k}(\mathbf{L}, \mathbf{S}, \mathbf{\Lambda}_k)$  using APG.
- 5:    $\mathbf{\Lambda}_{k+1} \leftarrow \mathbf{\Lambda}_k + \mu_k(\mathbf{L}_{k+1} + \mathbf{S}_{k+1} - \mathbf{Y})$ .
- 6:    $\mu_{k+1} \leftarrow \min\{\beta\mu_k, \mu_{\max}\}$ .
- 7: **end for**
- 8: **Output:**  $\mathbf{L}_\star \leftarrow \mathbf{L}_K, \mathbf{S}_\star \leftarrow \mathbf{S}_K$ .

# Optimization with Separable Structures

**Example:** Principal Component Pursuit

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \quad \text{subject to} \quad \mathbf{L} + \mathbf{S} = \mathbf{Y}. \quad (17)$$

A general two-term separable optimization program:

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}) \quad \text{subject to} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}, \quad (18)$$

where  $g$  and  $h$  are convex functions, and  $\mathbf{y} \in \text{range}([\mathbf{A} \mid \mathbf{B}])$ .

The Lagrangian  $\mathcal{L}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda})$  is:

$$\mathcal{L}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = g(\mathbf{x}) + h(\mathbf{z}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{y} \rangle. \quad (19)$$

# Optimization with Separable Structures

The augmented Lagrangian  $\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda})$  is:

$$\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = g(\mathbf{x}) + h(\mathbf{z}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} + \mathbf{Bz} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{y}\|_2^2. \quad (20)$$

**The alternating directions method of multipliers (ADMM)** conducts a simple, alternating iteration:

$$\mathbf{z}_{k+1} \in \arg \min_{\mathbf{z}} \mathcal{L}_\mu(\mathbf{x}_k, \mathbf{z}, \boldsymbol{\lambda}_k), \quad (21)$$

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}_{k+1}, \boldsymbol{\lambda}_k), \quad (22)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{Ax}_{k+1} + \mathbf{Bz}_{k+1} - \mathbf{y}). \quad (23)$$

This is also known as the *Gauss-Seidel iteration*.

**ADMM converges at a rate of  $O(1/k)$ .**

# ADMM for Principal Component Pursuit

$$\mathbf{PCP}: \min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \quad \text{subject to} \quad \mathbf{L} + \mathbf{S} = \mathbf{Y}. \quad (24)$$

The augmented Lagrangian is

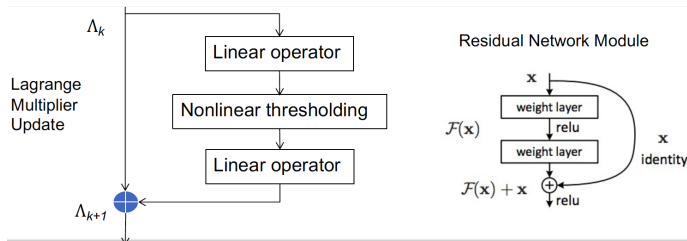
$$\mathcal{L}_\mu(\mathbf{L}, \mathbf{S}, \mathbf{\Lambda}) = \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \langle \mathbf{\Lambda}, \mathbf{L} + \mathbf{S} - \mathbf{Y} \rangle + \frac{\mu}{2} \|\mathbf{L} + \mathbf{S} - \mathbf{Y}\|_F^2. \quad (25)$$

$$\begin{aligned} \mathbf{L}_{k+1} &= \arg \min_{\mathbf{L}} \mathcal{L}_\mu(\mathbf{L}, \mathbf{S}_k, \mathbf{\Lambda}_k) \\ &= \arg \min_{\mathbf{L}} \|\mathbf{L}\|_* + \frac{\mu}{2} \|\mathbf{L} + \mathbf{S}_k - \mathbf{Y} + \mu^{-1} \mathbf{\Lambda}_k\|_F^2 + \varphi(\mathbf{S}_k, \mathbf{\Lambda}_k) \\ &= \text{prox}_{\mu^{-1} \|\cdot\|_*} [\mathbf{Y} - \mathbf{S}_k - \mu^{-1} \mathbf{\Lambda}_k]. \end{aligned} \quad (26)$$

$$\begin{aligned} \mathbf{S}_{k+1} &= \arg \min_{\mathbf{S}} \mathcal{L}_\mu(\mathbf{L}_{k+1}, \mathbf{S}, \mathbf{\Lambda}_k) \\ &= \arg \min_{\mathbf{S}} \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{S} + \mathbf{L}_{k+1} - \mathbf{Y} + \mu^{-1} \mathbf{\Lambda}_k\|_F^2 + \varphi(\mathbf{L}_{k+1}, \mathbf{\Lambda}_k) \\ &= \text{prox}_{\lambda \mu^{-1} \|\cdot\|_1} [\mathbf{Y} - \mathbf{L}_{k+1} - \mu^{-1} \mathbf{\Lambda}_k]. \end{aligned} \quad (27)$$

# ADMM Algorithm for PCP

- 1: **Problem:**  $\min_{L,S} \mathcal{L}_\mu(L, S, \Lambda)$ , given  $Y$ ,  $\lambda, \mu > 0$ .
- 2: **Input:**  $L_0, S_0, \Lambda_0 \in \mathbb{R}^{m \times n}$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:    $L_{k+1} \leftarrow \text{prox}_{\mu^{-1}\|\cdot\|_*} [Y - S_k - \mu^{-1}\Lambda_k]$ .
- 5:    $S_{k+1} \leftarrow \text{prox}_{\lambda\mu^{-1}\|\cdot\|_1} [Y - L_{k+1} - \mu^{-1}\Lambda_k]$ .
- 6:    $\Lambda_{k+1} \leftarrow \Lambda_k + \mu(L_{k+1} + S_{k+1} - Y)$ .
- 7: **end for**
- 8: **Output:**  $L_\star \leftarrow L_K; S_\star \leftarrow S_K$ .





## Multiple Separable Terms and Consensus Optimization

**Machine Learning:** Minimizing loss  $\sum_i L(\mathbf{y}_i, \mathbf{x})$  over samples  $\mathbf{y}_1, \dots, \mathbf{y}_p$ .  
We can partition the data to  $N$  batches, each on a machine:

$$\min_{\mathbf{x}} \sum_{j=1}^N f_j(\mathbf{x}) \quad \text{with} \quad f_j(\mathbf{x}) = \sum_{i \in I_j} L(\mathbf{y}_i, \mathbf{x}). \quad (28)$$

Convert to a consensus problem with separable variables:

$$\min_{\{\mathbf{x}_j\}} \sum_{j=1}^N f_j(\mathbf{x}_j) \quad \text{subject to} \quad \mathbf{x}_j = \mathbf{z}, \quad j = 1, \dots, N. \quad (29)$$

Augmented Lagrangian:

$$\mathcal{L}_{\mu}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = \sum_{j=1}^N f_j(\mathbf{x}_j) + \langle \boldsymbol{\lambda}_j, \mathbf{x}_j - \mathbf{z} \rangle + \frac{\mu}{2} \|\mathbf{x}_j - \mathbf{z}\|_2^2. \quad (30)$$

# Multiple Separable Terms and Consensus Optimization

Apply ADMM to the augmented Lagrangian  $\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda})$ :

$$\mathbf{x}_{j,k+1} = \arg \min_{\mathbf{x}_j} \left\{ f_j(\mathbf{x}_j) + \frac{\mu}{2} \left\| \mathbf{x}_j - \mathbf{z}_k + \frac{1}{\mu} \boldsymbol{\lambda}_{j,k} \right\|_2^2 \right\}, \text{ (parallel)} \quad (31)$$

$$\mathbf{z}_{k+1} = \frac{1}{N} \sum_{j=1}^N \left( \mathbf{x}_{j,k+1} + \frac{1}{\mu} \boldsymbol{\lambda}_{j,k} \right), \text{ (aggregate)} \quad (32)$$

$$\boldsymbol{\lambda}_{j,k+1} = \boldsymbol{\lambda}_{j,k} + \mu (\mathbf{x}_{j,k+1} - \mathbf{z}_{k+1}). \text{ (broadcast)} \quad (33)$$

**ADMM for ALM is well suited for distributed implementation!**

**Note:** there are many other variants to further improve efficiency and scalability: **accelerated, asynchronous, stochastic...** but convergence guarantee is not a picnic.

# Frank-Wolfe Algorithm

Optimizing a smooth, convex function over a *compact* convex set:

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{subject to} \quad \mathbf{x} \in \mathcal{C} \quad (34)$$

which has a finite diameter:

$$\text{diam}(\mathcal{C}) \doteq \max\{\|\mathbf{x} - \mathbf{x}'\|_2 \mid \mathbf{x}, \mathbf{x}' \in \mathcal{C}\}. \quad (35)$$

## Two examples:

- Sparse vector recovery:

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2, \quad \text{subject to} \quad \|\mathbf{x}\|_1 \leq \tau. \quad (36)$$

- Low-rank matrix completion:

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathcal{P}_{\Omega}[\mathbf{X}] - \mathbf{Y}\|_F^2, \quad \text{subject to} \quad \|\mathbf{X}\|_* \leq \tau. \quad (37)$$

# Franke-Wolfe Algorithm

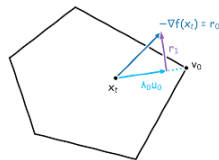
Find a point  $\mathbf{v}_k$  by solving a constrained optimization:

$$\mathbf{v}_k \in \arg \min_{\mathbf{v} \in \mathcal{C}} \langle \mathbf{v}, \nabla f(\mathbf{x}_k) \rangle. \quad (38)$$

We then set

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \gamma_k(\mathbf{v}_k - \mathbf{x}_k) = (1 - \gamma_k)\mathbf{x}_k + \gamma_k\mathbf{v}_k \in \mathcal{C}, \quad (39)$$

where  $\gamma_k \in (0, 1)$  is a specially chosen step size.



## Theorem (Convergence of Frank-Wolfe)

Let  $\mathbf{x}_0, \mathbf{x}_1, \dots$  denote the sequence of iterates generated by the Frank-Wolfe method, with step size  $\gamma_k = \frac{2}{k+2}$ . Then

$$f(\mathbf{x}_k) - f(\mathbf{x}_\star) \leq \frac{2L \text{diam}^2(\mathcal{C})}{k+2}. \quad (40)$$

# Frank-Wolfe for Matrix Completion

**Fact:** given a matrix  $G$  with SVD  $G = U\Sigma V^* = \sum_{i=1}^{n_1} u_i \sigma_i v_i$ , we have

$$V_* = -\tau \mathbf{u}_1 \mathbf{v}_1^* = \arg \min_V \langle V, G \rangle \quad \text{subject to} \quad \|V\|_* \leq \tau. \quad (41)$$

## Frank-Wolfe for Matrix Completion:

1: **Problem:** given  $Y = \mathcal{P}_\Omega[X_o + Z] \in \mathbb{R}^{n_1 \times n_2}$  and  $\Omega \subseteq [n_1] \times [n_2]$ ,

$$\min_X \frac{1}{2} \|\mathcal{P}_\Omega[X] - Y\|_F^2 \quad \text{subject to} \quad \|X\|_* \leq \tau.$$

2: **Input:**  $X_0 \in \mathbb{R}^{n_1 \times n_2}$  satisfying  $\|X_0\|_* \leq \tau$ .

3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**

4:    $(u_1, \sigma_1, v_1) \leftarrow \text{LeadSV}(\mathcal{P}_\Omega[X_k - Y])$  (power iteration).

5:    $V_k \leftarrow -\tau u_1 v_1^*$ .

6:    $X_{k+1} \leftarrow \frac{k}{k+2} X_k + \frac{2}{k+2} V_k$ .

7: **end for**

8: **Output:**  $X_* \leftarrow X_K$ .

# Franke-Wolfe for Noisy Sparse Recovery

1: **Problem:** given  $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathbf{z} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 \quad \text{subject to} \quad \|\mathbf{x}\|_1 \leq \tau.$$

2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$  satisfying  $\|\mathbf{x}_0\|_1 \leq \tau$ .

3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**

4:    $\mathbf{r}_k \leftarrow \mathbf{A}\mathbf{x}_k - \mathbf{y}$ .

5:    $i_k \leftarrow \arg \max_i |\mathbf{a}_i^* \mathbf{r}_k|$  (**matching pursuit**).

6:    $\sigma \leftarrow \text{sign}(\mathbf{a}_{i_k}^* \mathbf{r}_k)$ .

7:    $\mathbf{v}_k \leftarrow -\tau \sigma \mathbf{e}_{i_k}$ .

8:    $\mathbf{x}_{k+1} \leftarrow \frac{k}{k+2} \mathbf{x}_k + \frac{2}{k+2} \mathbf{v}_k$ .

9: **end for**

10: **Output:**  $\mathbf{x}_\star \leftarrow \mathbf{x}_K$ .

**Note:** Many greedy variants of the Franke-Wolfe algorithm: Matching Pursuit (MP), Orthogonal Matching Pursuit (OMP), Compressed Sampling Matching Pursuit (COSAMP), BLITZ, CELER etc.

## Other Ideas for Better Scalability

Typical optimization problem:  $\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m h_i(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n.$

**Complexity = per iteration cost  $\times$  # of iterations.**

- **Block Coordinate Descent** reduces dependency on the dimension  $n$ :

$$O(n) \rightarrow O(n^{1/2}).$$

- **Stochastic Gradient Descent** (with variance reduction) reduces dependency on sample size  $m$ :

$$O(m) \rightarrow O(m^{1/2}).$$

- **Acceleration Schemes** reduce the number of iterations  $k$ :

$$O(\epsilon^{-2}) \rightarrow O(\epsilon^{-1/2}).$$

# Assignments

- Reading: Section 8.4 - 8.6 of Chapter 8.
- Programming Homework #3.