Computational Principles for High-dim Data Analysis (Lecture Fourteen)

Yi Ma

EECS Department, UC Berkeley

October 14, 2021





EECS208, Fall 2021

Constrained & Scalable 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



• Challenge of Scale: scale algorithms to when n is very large.

Second order methods \implies First order methods... (2)

• Nonsmoothness: first order methods are slow for nonsmooth.

$$O(1/\sqrt{k}) \implies O(1/k) \implies O(1/k^2) \implies O(e^{-\alpha k})$$
 (3)

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

1

< □ > < □ > < □ > < □ > < □ > < □ >

Linear Equality Constrained Optimization

Problem:

$$\min_{\boldsymbol{x}} g(\boldsymbol{x}) \quad \text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}, \tag{4}$$

where

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

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

$$\hat{\boldsymbol{x}}(\mu) = \arg\min_{\boldsymbol{x}} |g(\boldsymbol{x}) + \frac{\mu}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|_2^2$$
 for a large μ . (5)

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

< ロ > < 同 > < 回 > < 回 > < 回 > <

Lagrange Multiplier Method A More Principled Approach:

Definition (The Lagrange Duality)

The Lagrangian function of the constrained problem (4):

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

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

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

Fact (credited to Lagrange): $\exists \lambda_{\star}$ such that the optimal solution $(x_{\star}, \lambda_{\star})$ is a saddle point of the Lagrangian:

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

< ロ > < 同 > < 回 > < 回 > < 回 > <

Dual Ascent Algorithm for the Lagrangian

Fact: If

$$oldsymbol{x}'(oldsymbol{\lambda}) = rg\min_{oldsymbol{x}} g(oldsymbol{x}) + \langle oldsymbol{\lambda}, oldsymbol{A}oldsymbol{x} - oldsymbol{y}
angle,$$

then $Ax'(\lambda) - y$ is a supergradient $\partial d(\lambda)$ of the concave dual $d(\lambda)$ at λ . (Actually this is true for the dual function of general constraints h(x) = 0: $d(\lambda) = \min_{x} g(x) + \lambda^{T} h(x)$. Why?)

A Natural Attempt to find the saddle point $(x_{\star}, \lambda_{\star})$ is via *dual ascent*:

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

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

- For certain problem classes, this converges to the optimal (x_\star, λ_\star) .
- However, unfortunately it fails for problems in our settings.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

An Example of Failure

Consider the basis pursuit problem:

$$\min_{\boldsymbol{x}} \|\boldsymbol{x}\|_1 \quad \text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}. \tag{11}$$

One can show that

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

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

The reason is g(x) here is not "strongly" convex to dominate the linear term $\langle \lambda, Ax \rangle$.



Augmented Lagrange Multiplier

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

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

Can be regarded as the Lagrangian for the penalized constrained problem

$$\min_{\boldsymbol{x}} \underbrace{g(\boldsymbol{x}) + \frac{\mu}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|_{2}^{2}}_{\text{strongly convex}} \quad \text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}, \quad (14)$$

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

(日) (四) (日) (日) (日)

Augmented Lagrange Multiplier

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

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

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

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

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

λ_{k+1} is always feasible, no bad behaviors!

イロト 不得 トイラト イラト 一日

Augmented Lagrange Multiplier

Augmented Lagrange Multipler (ALM)

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

Basic Iteration: set

$$\mathcal{L}_{\mu}(oldsymbol{x},oldsymbol{\lambda}) = g(oldsymbol{x}) + \langle oldsymbol{\lambda},oldsymbol{A}oldsymbol{x} - oldsymbol{y}
angle^2 + rac{\mu}{2} \left\|oldsymbol{A}oldsymbol{x} - oldsymbol{y}
ight\|_2^2.$$

Repeat:

$$egin{aligned} & m{x}_{k+1} \in rg\min_{m{x}} \ \mathcal{L}_{\mu}(m{x},m{\lambda}_k), \ & m{\lambda}_{k+1} = m{\lambda}_k + \mu \, (m{A}m{x}_{k+1} - m{y}). \end{aligned}$$

Convergence Guarantee:

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

(日)

ALM for Basis Pursuit

Augmented Lagrange Multipler (ALM) for BP

 Problem: min_x ||x||₁ subject to y = Ax, given y ∈ ℝ^m and A ∈ ℝ^{m×n}. The augmented Lagrangian is: L_μ(x, λ) = ||x||₁ + ⟨λ, Ax - y⟩ + μ/2 ||Ax - y||₂².
 Input: x₀ ∈ ℝⁿ, λ₀ ∈ ℝ^m, and β > 1.

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

4: $x_{k+1} \leftarrow \arg \min \mathcal{L}_{\mu_k}(x, \lambda_k)$ using APG.

5:
$$\boldsymbol{\lambda}_{k+1} \leftarrow \boldsymbol{\lambda}_k + \mu_k (\boldsymbol{A} \boldsymbol{x}_{k+1} - \boldsymbol{y}).$$

6:
$$\mu_{k+1} \leftarrow \min\{\beta \mu_k, \mu_{\max}\}.$$

- 7: end for
- 8: Output: $x_{\star} \leftarrow x_K$.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

ALM for Principal Component Pursuit

Augmented Lagrange Multipler (ALM) for PCP

1: **Problem:** $\min_{L,S} \|L\|_* + \lambda \|S\|_1$ subject to L + S = Y, given Y and $\lambda > 0$. The augmented Lagrangian:

$$\mathcal{L}_{\mu}(oldsymbol{L},oldsymbol{S},oldsymbol{\Lambda}) = \left\|oldsymbol{L}
ight\|_{*} + \lambda \left\|oldsymbol{S}
ight\|_{1} + \langleoldsymbol{\Lambda},oldsymbol{L} + oldsymbol{S} - oldsymbol{Y}
angle + rac{\mu}{2} \left\|oldsymbol{L} + oldsymbol{S} - oldsymbol{Y}
ight\|_{F}^{2}.$$

2: Input:
$$L_0, S_0, \Lambda_0 \in \mathbb{R}^{m \times n}$$
 and $\beta > 1$.
3: for $(k = 0, 1, 2, ..., K - 1)$ do
4: $\{L_{k+1}, S_{k+1}\} \leftarrow \arg \min \mathcal{L}_{\mu_k}(L, S, \Lambda_k)$ (APG? how? later...)
5: $\Lambda_{k+1} \leftarrow \Lambda_k + \mu_k(L_{k+1} + S_{k+1} - Y)$.
6: $\mu_{k+1} \leftarrow \min\{\beta \mu_k, \mu_{\max}\}$.
7: end for

8: Output: $L_{\star} \leftarrow L_K, S_{\star} \leftarrow S_K$.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

Optimization with Separable Structures

Example: Principal Component Pursuit

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

A general two-term separable optimization program:

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

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

The Lagrangian $\mathcal{L}(\boldsymbol{x}, \boldsymbol{z}, \boldsymbol{\lambda})$ has clear separable structures:

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

$$= \underbrace{g(\boldsymbol{x}) + \langle \boldsymbol{\lambda}, \boldsymbol{A} \boldsymbol{x} \rangle}_{\boldsymbol{\lambda}, \boldsymbol{\lambda}, \boldsymbol{\lambda}, \boldsymbol{\lambda}} + \underbrace{h(\boldsymbol{z}) + \langle \boldsymbol{\lambda}, \boldsymbol{B} \boldsymbol{z} \rangle}_{\boldsymbol{\lambda}, \boldsymbol{\lambda}, \boldsymbol{\lambda}$$

x dependent z

J

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Optimization with Separable Structures

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

$$\mathcal{L}_{\mu}(oldsymbol{x},oldsymbol{z},oldsymbol{\lambda}) = g(oldsymbol{x}) + h(oldsymbol{z}) + \langle oldsymbol{\lambda},oldsymbol{A}oldsymbol{x} + oldsymbol{B}oldsymbol{z} - oldsymbol{y}
angle + rac{\mu}{2} \|oldsymbol{A}oldsymbol{x} + oldsymbol{B}oldsymbol{z} - oldsymbol{y} \|_2^2.$$

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

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

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

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu \left(\boldsymbol{A} \boldsymbol{x}_{k+1} + \boldsymbol{B} \boldsymbol{z}_{k+1} - \boldsymbol{y} \right). \tag{23}$$

This is also known as the Gauss-Seidel iteration.

ADMM converges at a rate of O(1/k). (proof no picnic¹)

EECS208, Fall 2023

ADMM for Principal Component Pursuit

PCP: $\min_{L,S} \|L\|_* + \lambda \|S\|_1$ subject to L + S = Y. (24) The augmented Lagrangian is

$$\mathcal{L}_{\mu}(\boldsymbol{L},\boldsymbol{S},\boldsymbol{\Lambda}) = \|\boldsymbol{L}\|_{*} + \lambda \|\boldsymbol{S}\|_{1} + \langle \boldsymbol{\Lambda}, \boldsymbol{L} + \boldsymbol{S} - \boldsymbol{Y} \rangle + \frac{\mu}{2} \|\boldsymbol{L} + \boldsymbol{S} - \boldsymbol{Y}\|_{F}^{2}.$$
(25)
$$\boldsymbol{L}_{k+1} = \arg\min_{\boldsymbol{L}} \mathcal{L}_{\mu}(\boldsymbol{L}, \boldsymbol{S}_{k}, \boldsymbol{\Lambda}_{k})$$
$$= \arg\min_{\boldsymbol{L}} \|\boldsymbol{L}\|_{*} + \frac{\mu}{2} \|\boldsymbol{L} + \boldsymbol{S}_{k} - \boldsymbol{Y} + \mu^{-1}\boldsymbol{\Lambda}_{k}\|_{F}^{2} + \varphi(\boldsymbol{S}_{k}, \boldsymbol{\Lambda}_{k})$$
$$= \operatorname{prox}_{\mu^{-1}\|\cdot\|_{*}} [\boldsymbol{Y} - \boldsymbol{S}_{k} - \mu^{-1}\boldsymbol{\Lambda}_{k}].$$
(26)
$$\boldsymbol{S}_{k+1} = \arg\min_{\boldsymbol{S}} \mathcal{L}_{\mu}(\boldsymbol{L}_{k+1}, \boldsymbol{S}, \boldsymbol{\Lambda}_{k})$$
$$= \arg\min_{\boldsymbol{S}} \lambda \|\boldsymbol{S}\|_{1} + \frac{\mu}{2} \|\boldsymbol{S} + \boldsymbol{L}_{k+1} - \boldsymbol{Y} + \mu^{-1}\boldsymbol{\Lambda}_{k}\|_{F}^{2} + \varphi(\boldsymbol{L}_{k+1}, \boldsymbol{\Lambda}_{k})$$
$$= \operatorname{prox}_{\lambda\mu^{-1}\|\cdot\|_{1}} [\boldsymbol{Y} - \boldsymbol{L}_{k+1} - \mu^{-1}\boldsymbol{\Lambda}_{k}].$$
(27)

ADMM Algorithm for PCP

- 1: Problem: $\min_{\boldsymbol{L},\boldsymbol{S}} \mathcal{L}_{\mu}(\boldsymbol{L},\boldsymbol{S},\boldsymbol{\Lambda})$, given \boldsymbol{Y} , $\lambda, \mu > 0$. 2: Input: $\boldsymbol{L}_{0}, \boldsymbol{S}_{0}, \boldsymbol{\Lambda}_{0} \in \mathbb{R}^{m \times n}$. 3: for $(k = 0, 1, 2, \dots, K - 1)$ do 4: $\boldsymbol{L}_{k+1} \leftarrow \operatorname{prox}_{\mu^{-1} \parallel \cdot \parallel_{*}} [\boldsymbol{Y} - \boldsymbol{S}_{k} - \mu^{-1} \boldsymbol{\Lambda}_{k}]$. 5: $\boldsymbol{S}_{k+1} \leftarrow \operatorname{prox}_{\lambda \mu^{-1} \parallel \cdot \parallel_{*}} [\boldsymbol{Y} - \boldsymbol{L}_{k+1} - \mu^{-1} \boldsymbol{\Lambda}_{k}]$. 6: $\boldsymbol{\Lambda}_{k+1} \leftarrow \boldsymbol{\Lambda}_{k} + \mu(\boldsymbol{L}_{k+1} + \boldsymbol{S}_{k+1} - \boldsymbol{Y})$. 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 $\sum_i L(y_i, x)$ over samples y_1, \ldots, y_p . We can partition the data to N batches, each on a machine:

$$\min_{\boldsymbol{x}} \sum_{j=1}^{N} f_j(\boldsymbol{x}) \quad \text{with } f_j(\boldsymbol{x}) = \sum_{i \in \mathsf{I}_j} L(\boldsymbol{y}_i, \boldsymbol{x}). \tag{28}$$

Convert to a consensus problem with separable variables:

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

Augmented Lagrangian:

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

Multiple Separable Terms and Consensus Optimization

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

$$\begin{aligned} \boldsymbol{x}_{j,k+1} &= \arg\min_{\boldsymbol{x}_{j}} \left\{ f_{j}(\boldsymbol{x}_{j}) + \frac{\mu}{2} \| \boldsymbol{x}_{j} - \boldsymbol{z}_{k} + \frac{1}{\mu} \boldsymbol{\lambda}_{j,k} \|_{2}^{2} \right\}, \text{ (parallel) (31)} \\ \boldsymbol{z}_{k+1} &= \frac{1}{N} \sum_{j=1}^{N} \left(\boldsymbol{x}_{j,k+1} + \frac{1}{\mu} \boldsymbol{\lambda}_{j,k} \right), \text{ (aggregate)} \end{aligned}$$
(32)

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

< □ > < 同 > < 回 > < 回 > < 回 >

Frank-Wolfe Algorithm²

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

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}), \quad \text{subject to} \quad \boldsymbol{x} \in \mathsf{C} \tag{34}$$

which has a finite diameter:

$$\operatorname{diam}(\mathsf{C}) \doteq \max\{\|\boldsymbol{x} - \boldsymbol{x}'\|_2 \mid \boldsymbol{x}, \boldsymbol{x}' \in \mathsf{C}\}.$$
(35)

Two examples:

• Sparse vector recovery:

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

• Low-rank matrix completion:

$$\min_{\boldsymbol{X}} rac{1}{2} \left\| \mathcal{P}_{\Omega}[\boldsymbol{X}] - \boldsymbol{Y} \right\|_{F}^{2}, \hspace{0.1 in \ } ext{subject to} \hspace{0.1 in \ } \left\| \boldsymbol{X} \right\|_{*} \leq au.$$
 (37)

² "An algorithm for quadratic programming," M. Frank and P. Wolfe, Naval Research Logistics Quarterly, 1956.

Ma (EECS Department, UC Berkeley)

Franke-Wolfe Algorithm

Find a point v_k by solving a constrained optimization:

$$v_k \in \arg\min_{v \in \mathsf{C}} \langle v, \nabla f(x_k) \rangle.$$
 (complexity?) (38)

We then set



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

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

Theorem (Convergence of Frank-Wolfe)

Let x_0, x_1, \ldots denote the sequence of iterates generated by the Frank-Wolfe method, with step size $\gamma_k = \frac{2}{k+2}$. Then

$$f(\boldsymbol{x}_k) - f(\boldsymbol{x}_{\star}) \le \frac{2L \operatorname{diam}^2(\mathsf{C})}{k+2}.$$
(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

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

Frank-Wolfe for Matrix Completion:

1: **Problem:** given $\boldsymbol{Y} = \mathcal{P}_{\Omega}[\boldsymbol{X}_o + \boldsymbol{Z}] \in \mathbb{R}^{n_1 \times n_2}$ and $\Omega \subseteq [n_1] \times [n_2]$, $\min_{\boldsymbol{X}} \frac{1}{2} \|\mathcal{P}_{\Omega}[\boldsymbol{X}] - \boldsymbol{Y}\|_F^2 \quad \text{subject to} \quad \|\boldsymbol{X}\|_* \leq \tau.$

2: Input: $X_0 \in \mathbb{R}^{n_1 \times n_2}$ satisfying $||X_0||_* \leq \tau$. 3: for (k = 0, 1, 2, ..., 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: $X_{\star} \leftarrow X_K$.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Franke-Wolfe for Noisy Sparse Recovery

1: **Problem:** given $y = Ax_0 + z \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times n}$,

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

2: Input:
$$x_0 \in \mathbb{R}^n$$
 satisfying $\|x_0\|_1 \leq \tau$.
3: for $(k = 0, 1, 2, \dots, K - 1)$ do

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

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

- $i_k \leftarrow \arg \max_i |\boldsymbol{a}_i^* \boldsymbol{r}_k|$ (matching pursuit). 5:
- $\sigma \leftarrow \operatorname{sign}\left(\boldsymbol{a}_{i_{1}}^{*} \boldsymbol{r}_{k}\right).$ 6:

7:
$$v_k \leftarrow -\tau \sigma e_{i_k}$$
.

8:
$$x_{k+1} \leftarrow rac{k}{k+2} x_k + rac{2}{k+2} v_k$$
.

9: end for

10: Output: $x_{\star} \leftarrow x_{K}$.

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

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Other Ideas for Better Scalability

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

Complexity = per iteration cost \times # of iterations.

• Block Coordinate Descent reduces dependency on the dimension n:

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

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

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

• Acceleration Schemes reduce the number of iterations k:

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

Nonconvex programs are a different story... (later, Chapter 9).

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Assignments

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

э

-

< /□ > < Ξ