

Alternating Direction Method of Multipliers

a talk by

Vinícius and Prof. Daniel Palomar

The Hong Kong University of Science and Technology

ELEC5470/IEDA6100A - Convex Optimization

Contents

1. Introduction

Optimization algorithms, motivation

2. Alternating Direction Method of Multipliers

The basics

3. Practical Examples

Robust PCA and Graphical Lasso

Why use optimization algorithms?

Motivations

- ❖ large-scale optimization
 - ❖ machine learning/statistics with huge datasets
 - ❖ computer vision
- ❖ decentralized optimization
 - ❖ entities/agents/threads coordinate to solve a large problem by passing small messages

Optimization Algorithms

- ❖ Gradient Descent
- ❖ Newton
- ❖ Interior Point Methods (IPM)
- ❖ Block Coordinate Descent (BCD)
- ❖ Majorization-Minimization (MM)
- ❖ Block Majorization-Minimization (BMM)
- ❖ Successive Convex Approximation (SCA)

Optimization Algorithms

- ❖ Gradient Descent
- ❖ Newton
- ❖ Interior Point Methods (IPM)
- ❖ Block Coordinate Descent (BCD)
- ❖ Majorization-Minimization (MM)
- ❖ Block Majorization-Minimization (BMM)
- ❖ Successive Convex Approximation (SCA)
- ❖ ...

Optimization Algorithms

- ❖ Gradient Descent
- ❖ Newton
- ❖ Interior Point Methods (IPM)
- ❖ Block Coordinate Descent (BCD)
- ❖ Majorization-Minimization (MM)
- ❖ Block Majorization-Minimization (BMM)
- ❖ Successive Convex Approximation (SCA)
- ❖ ...
- ❖ **Alternating Direction Method of Multipliers (ADMM)**

Reference

- ❖ Boyd *et al.* **Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers.** *Foundations and Trends in Machine Learning*. 2010.
- ❖ available online for **free**: https://web.stanford.edu/~boyd/papers/pdf/admm_distr_stats.pdf
- ❖ citations: 13519¹
- ❖ Boyd's presentation: https://web.stanford.edu/class/ee364b/lectures/admm_slides.pdf
- ❖ Yuxin Chen's Princeton lecture notes ELE 522: Large-Scale Optimization for Data Science

¹as of Nov. 24th 2020

Dual Problem

- convex equality constrained optimization problem

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{Ax} = \mathbf{b}\end{array}$$

- Lagrangian: $L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top (\mathbf{Ax} - \mathbf{b})$
- dual function: $g(\mathbf{y}) = \inf_{\mathbf{x}} L(\mathbf{x}, \mathbf{y})$
- dual problem: maximize $g(\mathbf{y})$
- recover: $\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} L(\mathbf{x}, \mathbf{y}^*)$

Dual Ascent

Dual Ascent

- gradient method for dual problem: $\mathbf{y}^{k+1} = \mathbf{y}^k + \rho^k \nabla g(\mathbf{y}^k)$
- $\nabla g(\mathbf{y}^k) = \mathbf{A}\mathbf{x}^{k+1} - \mathbf{b}$, where $\mathbf{x}^{k+1} = \arg \min_{\mathbf{x}} L(\mathbf{x}, \mathbf{y}^k)$
- dual ascent method is

$$\mathbf{x}^{k+1} := \arg \min_{\mathbf{x}} L(\mathbf{x}, \mathbf{y}^k)$$

$$\mathbf{y}^{k+1} := \mathbf{y}^k + \rho^k (\mathbf{A}\mathbf{x}^{k+1} - \mathbf{b})$$

Dual Ascent

- gradient method for dual problem: $\mathbf{y}^{k+1} = \mathbf{y}^k + \rho^k \nabla g(\mathbf{y}^k)$
- $\nabla g(\mathbf{y}^k) = \mathbf{A}\mathbf{x}^{k+1} - \mathbf{b}$, where $\mathbf{x}^{k+1} = \arg \min_{\mathbf{x}} L(\mathbf{x}, \mathbf{y}^k)$
- dual ascent method is

$$\mathbf{x}^{k+1} := \arg \min_{\mathbf{x}} L(\mathbf{x}, \mathbf{y}^k)$$

$$\mathbf{y}^{k+1} := \mathbf{y}^k + \rho^k (\mathbf{A}\mathbf{x}^{k+1} - \mathbf{b})$$

- why?

Dual Decomposition

Dual Decomposition

- ❖ suppose f is separable:

$$f(\mathbf{x}) = f_1(x_1) + \cdots + f_n(x_n), \mathbf{x} = (x_1, \dots, x_n)$$

- ❖ then the Lagrangian is separable in \mathbf{x} :

$$L_i(x_i, \mathbf{y}) = f_i(x_i) + \mathbf{y}^\top \mathbf{a}_{*,i} x_i$$

- ❖ \mathbf{x} -minimization splits into n separate minimizations

$$x_i^{k+1} := \arg \min_{x_i} L_i(x_i, \mathbf{y}^k), i = 1, \dots, n$$

which can be done in parallel and $\mathbf{y}^{k+1} = \mathbf{y}^k + \alpha^k (\sum_{i=1}^n \mathbf{a}_{*,i} x_i^{k+1} - \mathbf{b})$

Optimization Problem

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} && f(\mathbf{x}) + g(\mathbf{z}) \\ & \text{subject to} && \mathbf{Ax} + \mathbf{Bz} = \mathbf{c} \end{aligned} \tag{1}$$

- ❖ variables: $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{z} \in \mathbb{R}^m$
- ❖ parameters: $\mathbf{A} \in \mathbb{R}^{p \times n}$, $\mathbf{B} \in \mathbb{R}^{p \times m}$, and $\mathbf{c} \in \mathbb{R}^p$
- ❖ optimal value: $p^* = \inf_{\mathbf{x}, \mathbf{z}} \{f(\mathbf{x}) + g(\mathbf{z}) : \mathbf{Ax} + \mathbf{Bz} = \mathbf{c}\}$

Augmented Lagrangian Method

Augmented Lagrangian Method

- Augmented Lagrangian:

$$L_{\rho}(\mathbf{x}, \mathbf{z}, \mathbf{y}) = \underbrace{f(\mathbf{x}) + g(\mathbf{z}) + \langle \mathbf{y}, \mathbf{Ax} + \mathbf{Bz} - \mathbf{c} \rangle}_{\text{Lagrangian}} + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}\|_{\text{F}}^2$$

- ALM consists of the iterations:

$$\mathbf{x}^{k+1}, \mathbf{z}^{k+1} := \arg \min_{\mathbf{x}, \mathbf{z}} L_{\rho}(\mathbf{x}, \mathbf{z}, \mathbf{y}^k) \text{ // primal update}$$

$$\mathbf{y}^{k+1} := \mathbf{y}^k + \rho (\mathbf{Ax}^{k+1} + \mathbf{Bz}^{k+1} - \mathbf{c}) \text{ // dual update}$$

- $\rho > 0$ is a penalty hyperparameter

Issues with Augmented Lagrangian Method

- ❖ the primal step is often expensive to solve – as expensive as solving the original problem
- ❖ minimization of \mathbf{x} and \mathbf{z} has to be done jointly

Alternating Direction Method of Multipliers

Alternating Direction Method of Multipliers

- Augmented Lagrangian:

$$L_{\rho}(\mathbf{x}, \mathbf{z}, \mathbf{y}) = \underbrace{f(\mathbf{x}) + g(\mathbf{z}) + \langle \mathbf{y}, \mathbf{Ax} + \mathbf{Bz} - \mathbf{c} \rangle}_{\text{Lagrangian}} + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}\|_{\text{F}}^2$$

- ADMM consists of the iterations:

$$\mathbf{x}^{k+1} := \arg \min_{\mathbf{x}} L_{\rho}(\mathbf{x}, \mathbf{z}^k, \mathbf{y}^k)$$

$$\mathbf{z}^{k+1} := \arg \min_{\mathbf{z}} L_{\rho}(\mathbf{x}^{k+1}, \mathbf{z}, \mathbf{y}^k)$$

$$\mathbf{y}^{k+1} := \mathbf{y}^k + \rho (\mathbf{Ax}^{k+1} + \mathbf{Bz}^{k+1} - \mathbf{c})$$

- $\rho > 0$ is a penalty hyperparameter

Convergence and Stopping Criteria

- ❖ assume (very little!)
 - ❖ f, g are convex, closed, proper
 - ❖ L_0 has a saddle point
- ❖ then ADMM converges:
 - ❖ iterates approach feasibility: $\mathbf{Ax}^k + \mathbf{Bz}^k - \mathbf{c} \rightarrow \mathbf{0}$
 - ❖ objective approaches optimal value: $f(\mathbf{x}^k) + g(\mathbf{z}^k) \rightarrow p^*$
- ❖ false (in general) statements: \mathbf{x} converges, \mathbf{z} converges
- ❖ true statement: \mathbf{y} converges
- ❖ what matters: residual is small and near optimality in objective value

Convergence of ADMM in Practice

- ❖ ADMM is often slow to converge to high accuracy
- ❖ ADMM often converges to moderate accuracy within a few dozens of iterations, which is often sufficient for most practical purposes

Practical Examples

Robust PCA (Candes et al. '08)

- ❖ We would like to model a data matrix \mathbf{M} as low-rank plus sparse components:

$$\begin{aligned} & \underset{\mathbf{L}, \mathbf{S}}{\text{minimize}} && \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \\ & \text{subject to} && \mathbf{L} + \mathbf{S} = \mathbf{M} \end{aligned}$$

- ❖ where $\|\mathbf{L}\|_* := \sum_{i=1}^n \sigma_i(\mathbf{L})$ is the nuclear norm
- ❖ and $\|\mathbf{S}\|_1 := \sum_{i,j} |S_{ij}|$ is the entrywise ℓ_1 -norm

Robust PCA via ADMM

ADMM for solving robust PCA:

$$\mathbf{L}^{k+1} = \arg \min_{\mathbf{L}} \|\mathbf{L}\|_* + \text{tr}(\mathbf{Y}^{k\top} \mathbf{L}) + \frac{\rho}{2} \|\mathbf{L} + \mathbf{S}^k - \mathbf{M}\|_{\text{F}}^2$$

$$\mathbf{S}^{k+1} = \arg \min_{\mathbf{S}} \lambda \|\mathbf{S}\|_1 + \text{tr}(\mathbf{Y}^{k\top} \mathbf{S}) + \frac{\rho}{2} \|\mathbf{L}^{k+1} + \mathbf{S} - \mathbf{M}\|_{\text{F}}^2$$

$$\mathbf{Y}^{k+1} = \mathbf{Y}^k + \rho (\mathbf{L}^{k+1} + \mathbf{S}^{k+1} - \mathbf{M})$$

Robust PCA via ADMM

$$\begin{aligned}\mathbf{L}^{k+1} &= \text{SVT}_{\rho^{-1}} \left(\mathbf{M} - \mathbf{s}^k - \frac{1}{\rho} \mathbf{Y}^k \right) \\ \mathbf{s}^{k+1} &= \text{ST}_{\lambda \rho^{-1}} \left(\mathbf{M} - \mathbf{L}^{k+1} - \frac{1}{\rho} \mathbf{Y}^k \right) \\ \mathbf{Y}^{k+1} &= \mathbf{Y}^k + \rho \left(\mathbf{L}^{k+1} + \mathbf{s}^{k+1} - \mathbf{M} \right),\end{aligned}$$

where for any \mathbf{X} with SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, $\mathbf{\Sigma} = \text{diag}(\{\sigma_i\})$, we have

$$\text{SVT}_{\tau}(\mathbf{X}) = \mathbf{U} \text{diag}(\{(\sigma_i - \tau)^+\}) \mathbf{V}^\top$$

and

$$(\text{ST}_{\tau}(\mathbf{X}))_{ij} = \begin{cases} X_{ij} - \tau, & \text{if } X_{ij} > \tau, \\ 0, & \text{if } |X_{ij}| \leq \tau, \\ X_{ij} + \tau, & \text{if } X_{ij} < -\tau \end{cases}$$

Graphical Lasso

Precision matrix estimation from Gaussian samples:

$$\begin{aligned} & \underset{\Theta}{\text{minimize}} && \underbrace{-\log \det \Theta + \langle \Theta, \mathbf{S} \rangle}_{\text{neg. log likelihood}} + \lambda \|\Theta\|_1 \\ & \text{subject to} && \Theta \succ \mathbf{0} \end{aligned}$$

Or equivalently, using a slack variable $\Psi = \Theta$

$$\begin{aligned} & \underset{\Theta, \Psi}{\text{minimize}} && \underbrace{-\log \det \Theta + \langle \Theta, \mathbf{S} \rangle}_{\text{neg. log likelihood}} + \lambda \|\Psi\|_1 \\ & \text{subject to} && \Theta \succ \mathbf{0}, \Theta = \Psi \end{aligned}$$

Graphical Lasso via ADMM

$$\Theta^{k+1} = \arg \min_{\Theta \succ \mathbf{0}} -\log \det \Theta + \langle \Theta, \mathbf{S} + \mathbf{Y}^k \rangle + \frac{\rho}{2} \left\| \Theta - \Psi^k \right\|_F^2$$

$$\Psi^{k+1} = \arg \min_{\Psi} \lambda \left\| \Psi \right\|_1 - \langle \Psi, \mathbf{Y}^k \rangle + \frac{\rho}{2} \left\| \Theta^k - \Psi \right\|_F^2$$

$$\mathbf{Y}^{k+1} = \mathbf{Y}^k + \rho \left(\Theta^{k+1} - \Psi^{k+1} \right)$$

Graphical Lasso via ADMM

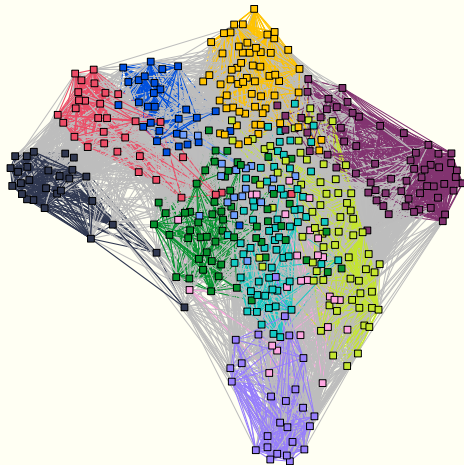
$$\Theta^{k+1} = \mathcal{F}_\rho \left(\Psi^k - \frac{1}{\rho} (\mathbf{Y}^k + \mathbf{s}) \right)$$

$$\Psi^{k+1} = \text{ST}_{\lambda\rho^{-1}} \left(\Theta^{k+1} + \frac{1}{\rho} \mathbf{Y}^k \right)$$

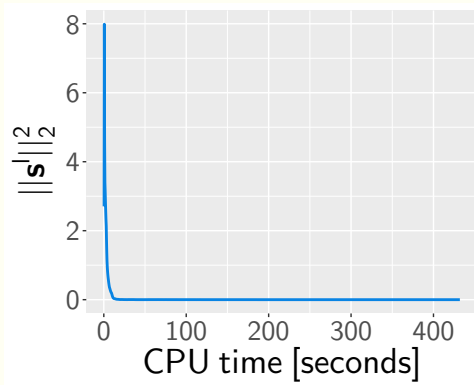
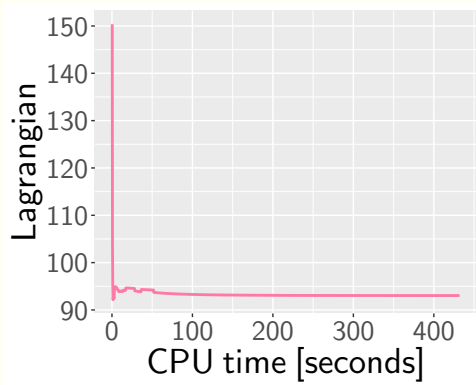
$$\mathbf{Y}^{k+1} = \mathbf{Y}^k + \rho (\Theta^{k+1} - \Psi^{k+1})$$

❖ where $\mathcal{F}_\rho(\mathbf{X}) := \frac{1}{2} \mathbf{U} \text{diag} \left(\left\{ \lambda_i + \sqrt{\lambda_i^2 + \frac{4}{\rho}} \right\} \right) \mathbf{U}^\top$, for $\mathbf{X} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$.

Network of stocks via Graphical Lasso



Network of stocks via Graphical Lasso



Conclusion

- ❖ ADMM is a versatile/flexible optimization framework
- ❖ may not be the best for a specific case, but often performs well in practice
- ❖ convergence often needs to be proved in a case-by-case scenario

Questions?