

Alternating direction method of multipliers

Ye Shi

ShanghaiTech University

Outline

- Augmented Lagrangian method
- Alternating direction method of multipliers



Two-block problem

$$\begin{array}{ll} \text{minimize}_{\boldsymbol{x}, \boldsymbol{z}} & F(\boldsymbol{x}, \boldsymbol{z}) := f_1(\boldsymbol{x}) + f_2(\boldsymbol{z}) \\ \text{subject to} & \boldsymbol{A}\boldsymbol{x} + \boldsymbol{B}\boldsymbol{z} = \boldsymbol{b} \end{array}$$

where f_1 and f_2 are both convex

- this can also be solved via Douglas-Rachford splitting
- we will introduce another paradigm for solving this problem

Augmented Lagrangian method

Dual problem

$$\begin{array}{ll} \text{minimize}_{x,z} & f_1(x) + f_2(z) \checkmark \\ \text{subject to} & \lambda(Ax + Bz = b) \end{array} \quad \left. \vphantom{\begin{array}{l} \text{minimize}_{x,z} \\ \text{subject to} \end{array}} \right\}$$

\Updownarrow

$$\text{maximize}_{\lambda} \quad \underbrace{\min_{x,z} f_1(x) + f_2(z) + \langle \lambda, Ax + Bz - b \rangle}_{=: \mathcal{L}(x,z,\lambda) \text{ (Lagrangian)}} \quad \checkmark$$

\Updownarrow

$$\text{maximize}_{\lambda} \quad -f_1^*(-A^\top \lambda) - f_2^*(-B^\top \lambda) - \langle \lambda, b \rangle \quad \checkmark \checkmark$$

\Updownarrow

$$\text{minimize}_{\lambda} \quad f_1^*(-A^\top \lambda) + f_2^*(-B^\top \lambda) + \langle \lambda, b \rangle \quad \checkmark$$

Augmented Lagrangian method ✓

$$\text{minimize}_{\lambda} \quad \underbrace{f_1^*(-A^\top \lambda)} + \underbrace{f_2^*(-B^\top \lambda)} + \langle \lambda, b \rangle$$

The proximal point method for solving this dual problem:

$$\lambda^{t+1} = \arg \min_{\lambda} \left\{ \underbrace{f_1^*(-A^\top \lambda)} + \underbrace{f_2^*(-B^\top \lambda)} + \langle \lambda, b \rangle + \frac{1}{2\rho} \|\lambda - \lambda^t\|_2^2 \right\} \quad \delta$$



As it turns out, this is equivalent to the **augmented Lagrangian method** (or the method of multipliers)

$$\begin{aligned} (x^{t+1}, z^{t+1}) &= \arg \min_{x, z} \left\{ \underbrace{f_1(x)} + \underbrace{f_2(z)} + \frac{\rho}{2} \left\| Ax + Bz - b + \frac{1}{\rho} \lambda^t \right\|_2^2 \right\} \quad \checkmark \\ \lambda^{t+1} &= \lambda^t + \rho \underbrace{(Ax^{t+1} + Bz^{t+1} - b)}_{\delta} \quad \checkmark \end{aligned} \quad (10.1)$$

λ^*

$$\begin{cases} \min & f(x) \\ \text{s.t.} & h(x) = 0 \end{cases} \quad \underline{\text{optimal}}$$

$$\left\{ \min \quad \underbrace{L(x, \lambda) = f(x) + \lambda h(x)}_{\text{①}} \quad \right\}$$

$$\text{loss:} \quad \text{Cross-Entropy} + \lambda \text{ Regularization}$$

Justification of (10.1)

$$\lambda^{t+1} = \arg \min_{\lambda} \left\{ f_1^*(-A^\top \lambda) + f_2^*(-B^\top \lambda) + \langle \lambda, b \rangle + \frac{1}{2\rho} \|\lambda - \lambda^t\|_2^2 \right\} \quad \checkmark$$

\Updownarrow optimality condition

$$0 \in -A \partial f_1^*(-A^\top \lambda^{t+1}) - \underbrace{B \partial f_2^*(-B^\top \lambda^{t+1}) + b + \frac{1}{\rho} (\lambda^{t+1} - \lambda^t)}_{\text{subdifferential}}$$

\Updownarrow

$$\lambda^{t+1} = \lambda^t + \rho (Ax^{t+1} + Bz^{t+1} - b) \quad \checkmark$$

where (check: use the conjugate subgradient theorem)

$$\underbrace{x^{t+1}} := \arg \min_x \{ \langle \underbrace{A^\top \lambda^{t+1}}, x \rangle + f_1(x) \} \quad \checkmark -$$

$$z^{t+1} := \arg \min_z \{ \langle B^\top \lambda^{t+1}, z \rangle + f_2(z) \} \quad \checkmark -$$

$$0 \in A^\top \lambda^{t+1} + \partial f_1(x^{t+1}) \Rightarrow -A^\top \lambda^{t+1} \in \partial f_1(x^{t+1}) \Rightarrow x^{t+1} \in \partial f_1^*(-A^\top \lambda^{t+1})$$

Justification of (10.1)

 \Updownarrow

$$\begin{aligned}\mathbf{x}^{t+1} &:= \arg \min_{\mathbf{x}} \left\{ \langle \mathbf{A}^\top [\boldsymbol{\lambda}^t + \rho (\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^{t+1} - \mathbf{b})], \mathbf{x} \rangle + f_1(\mathbf{x}) \right\} \\ \mathbf{z}^{t+1} &:= \arg \min_{\mathbf{z}} \left\{ \langle \mathbf{B}^\top [\boldsymbol{\lambda}^t + \rho (\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^{t+1} - \mathbf{b})], \mathbf{z} \rangle + f_2(\mathbf{z}) \right\}\end{aligned}$$

 \Updownarrow

$$\begin{aligned}\mathbf{0} &\in \mathbf{A}^\top [\boldsymbol{\lambda}^t + \rho (\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^{t+1} - \mathbf{b})] + \partial f_1(\mathbf{x}^{t+1}) \\ \mathbf{0} &\in \mathbf{B}^\top [\boldsymbol{\lambda}^t + \rho (\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^{t+1} - \mathbf{b})] + \partial f_2(\mathbf{z}^{t+1})\end{aligned}$$

 \Updownarrow

$$(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}) = \arg \min_{\mathbf{x}, \mathbf{z}} \left\{ \underbrace{f_1(\mathbf{x}) + f_2(\mathbf{z}) + \frac{\rho}{2} \left\| \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right\|_2^2}_{\text{red wavy line}} \right\}$$

Augmented Lagrangian method (ALM)

$$\underbrace{(x^{t+1}, z^{t+1})}_{-} = \arg \min_{\underbrace{x, z}_{-}} \left\{ \underbrace{f_1(x)}_{-} + \underbrace{f_2(z)}_{-} + \frac{\rho}{2} \left\| Ax + Bz - b + \frac{1}{\rho} \lambda^t \right\|_2^2 \right\}$$

(primal step)

$$\underbrace{\lambda^{t+1}}_{-} = \lambda^t + \rho (Ax^{t+1} + Bz^{t+1} - b)$$

(dual step)

where $\rho > 0$ is penalty parameter

ALM aims to solve the following problem by alternating between primal and dual updates

$$\text{maximize}_{\lambda} \underbrace{\max_{x, z} f_1(x) + f_2(z) + \rho \langle Ax + Bz - b, \lambda \rangle + \frac{\rho}{2} \left\| Ax + Bz - b + \frac{1}{\rho} \lambda \right\|_2^2}_{\mathcal{L}_{\rho}(x, z, \lambda): \text{ augmented Lagrangian}}$$

Issues of augmented Lagrangian method

$$(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}) = \arg \min_{\mathbf{x}, \mathbf{z}} \left\{ \underbrace{f_1(\mathbf{x})}_{\text{minimize}} + \underbrace{f_2(\mathbf{z})}_{\text{minimize}} + \frac{\rho}{2} \left\| \mathbf{Ax} + \mathbf{Bz} - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right\|_2^2 \right\} \quad \checkmark$$

- the primal update step is often expensive — as expensive as solving the original problem
- minimization of \mathbf{x} and \mathbf{z} cannot be carried out separately \checkmark

Alternating direction method of multipliers

Alternating direction method of multipliers

Rather than computing exact primal estimate for ALM, we might minimize x and z sequentially via alternating minimization

$$\left. \begin{aligned} \underset{x}{x}^{t+1} &= \arg \min_x \left\{ f_1(x) + \frac{\rho}{2} \left\| Ax + Bz^t - b + \frac{1}{\rho} \lambda^t \right\|_2^2 \right\} \\ \underset{z}{z}^{t+1} &= \arg \min_z \left\{ f_2(z) + \frac{\rho}{2} \left\| Ax^{t+1} + Bz - b + \frac{1}{\rho} \lambda^t \right\|_2^2 \right\} \\ \lambda^{t+1} &= \lambda^t + \rho (Ax^{t+1} + Bz^{t+1} - b) \end{aligned} \right\}$$

— called the *alternating direction method of multipliers (ADMM)*

$$\|x^{t+1} - x^*\| \leq \dots \|f(x^{t+1}) - f(x^*)\|$$

Alternating direction method of multipliers

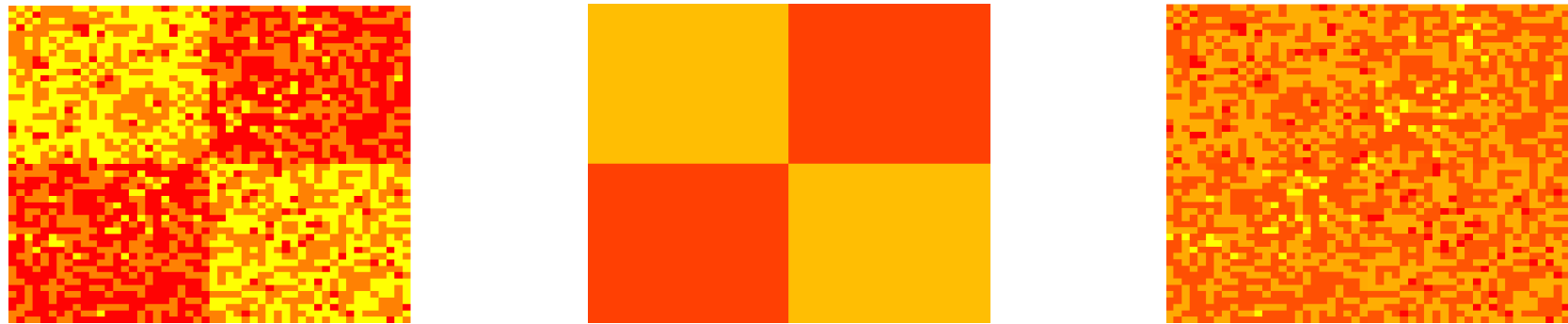
$$\mathbf{x}^{t+1} = \arg \min_{\mathbf{x}} \left\{ \underbrace{f_1(\mathbf{x})}_{\text{red wavy}} + \frac{\rho}{2} \left\| \mathbf{Ax} + \mathbf{Bz}^t - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right\|_2^2 \right\} \quad \checkmark$$

$$\mathbf{z}^{t+1} = \arg \min_{\mathbf{z}} \left\{ f_2(\mathbf{z}) + \frac{\rho}{2} \left\| \mathbf{Ax}^{t+1} + \mathbf{Bz} - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right\|_2^2 \right\} \quad \checkmark \checkmark$$

$$\boldsymbol{\lambda}^{t+1} = \boldsymbol{\lambda}^t + \rho(\mathbf{Ax}^{t+1} + \mathbf{Bz}^{t+1} - \mathbf{b})$$

- ρ controls relative priority between primal and dual convergence
- useful if updating \mathbf{x}^t and updating \mathbf{z}^t are both inexpensive
- blend the benefits of dual decomposition and augmented Lagrangian method
- the roles of \mathbf{x} and \mathbf{z} are *almost* symmetric, but not quite

Example: robust PCA


$$\underbrace{M}_{\text{observed}} = \underbrace{L}_{\text{low-rank}} + \underbrace{S}_{\text{sparse}}$$

Suppose we observe M , which is the superposition of a low-rank component L and sparse outliers S

Can we hope to disentangle L and S ?

Example: robust PCA

One way to solve it is via convex programming (Candes et al. '08)

$$\begin{aligned} \underset{L, S}{\text{minimize}} \quad & \overbrace{\|L\|_*}^{f_1(x)} + \lambda \overbrace{\|S\|_1}^{f_2(z)} \\ \text{s.t.} \quad & \underbrace{L + S = M} \end{aligned} \quad (10.2)$$

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

$$\min_L \quad \underbrace{\|L\|_* + \frac{\rho}{2} \|L + S^t - M\|_F^2 + \frac{1}{\rho} \lambda^t}_{\text{ADMM}}$$

Example: robust PCA

ADMM for solving (10.2):

$$\begin{aligned}\mathbf{L}^{t+1} &= \arg \min_{\mathbf{L}} \left\{ \underbrace{\|\mathbf{L}\|_* + \frac{\rho}{2} \|\mathbf{L} + \mathbf{S}^t - \mathbf{M} + \frac{1}{\rho} \mathbf{\Lambda}^t\|_{\text{F}}^2}_{\Delta} \right\} \smile \\ \mathbf{S}^{t+1} &= \arg \min_{\mathbf{S}} \left\{ \lambda \|\mathbf{S}\|_1 + \frac{\rho}{2} \|\mathbf{L}^{t+1} + \mathbf{S} - \mathbf{M} + \frac{1}{\rho} \mathbf{\Lambda}^t\|_{\text{F}}^2 \right\} \smile \\ \mathbf{\Lambda}^{t+1} &= \mathbf{\Lambda}^t + \rho(\mathbf{L}^{t+1} + \mathbf{S}^{t+1} - \mathbf{M})\end{aligned}$$

Example: robust PCA

This is equivalent to

$$\mathbf{L}^{t+1} = \text{SVT}_{\rho^{-1}} \left(\mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t \right) \quad (\text{singular value thresholding}) \quad \checkmark$$

$$\mathbf{S}^{t+1} = \text{ST}_{\lambda \rho^{-1}} \left(\mathbf{M} - \mathbf{L}^{t+1} - \frac{1}{\rho} \mathbf{\Lambda}^t \right) \quad (\text{soft thresholding}) \quad \checkmark$$

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

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

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

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

$$\text{prox}_{\|\cdot\|_*}(x) = \arg \min_{\underline{z} \in \mathbb{R}^{m \times n}} \frac{1}{2} \|\underline{x} - \underline{z}\|_F^2 + \lambda \|\underline{z}\|_* \quad \checkmark$$

$$\underline{X} = \underline{U} \underline{\Sigma} \underline{V}^T \quad (\text{SVD decomposition})$$

$$\underline{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_{\min(m,n)}) \quad \checkmark$$

$$\arg \min_{\underline{z} \in \mathbb{R}^{m \times n}} \frac{1}{2} \|\underline{U} \underline{\Sigma} \underline{V}^T - \underline{z}\|_F^2 + \lambda \|\underline{z}\|_*$$

$$= \arg \min_{\underline{z} \in \mathbb{R}^{m \times n}} \frac{1}{2} \|\underline{\Sigma} - \underline{U}^T \underline{z} \underline{V}\|_F^2 + \lambda \|\underline{z}\|_*$$

$$\underline{\tilde{z}} := \underline{U}^T \underline{z} \underline{V} \Rightarrow \underline{z} = \underline{U} \underline{\tilde{z}} \underline{V}^T$$

$$\text{prox}_{\|\cdot\|_*}(x) = \underline{U} \left(\arg \min_{\underline{\tilde{z}}} \frac{1}{2} \|\underline{\Sigma} - \underline{\tilde{z}}\|_F^2 + \lambda \|\underline{U} \underline{\tilde{z}} \underline{V}^T\|_* \right) \underline{V}^T$$

$$\|\underline{\tilde{z}}\|_* = \|\underline{z}\|_* \quad , \quad \underline{\tilde{z}} = \underline{U}_{\tilde{z}} \underline{\Sigma}_{\tilde{z}} \underline{V}_{\tilde{z}}^T \quad ,$$

$$\underline{z} = \underline{U} \underline{\tilde{z}} \underline{V}^T = (\underline{U} \underline{U}_{\tilde{z}}) \underline{\Sigma}_{\tilde{z}} (\underline{V}_{\tilde{z}}^T \underline{V}^T)$$

$$\text{prox}_{\|\cdot\|_*}(x) = U \left(\underset{\tilde{Z}}{\text{argmin}} \frac{1}{2} \|\varepsilon - \tilde{Z}\|_F^2 + \lambda \|\tilde{Z}\|_* \right) U^T$$

$\text{prox}_{\|\cdot\|_*}(\varepsilon)$

\tilde{Z} is diagonal to ensure the minimization

$$\underset{\tilde{Z}}{\text{argmin}}_{\|\cdot\|_*} \left[\sum_i \left[\frac{1}{2} (\tilde{z}_{ii} - \varepsilon_i)^2 + \lambda \tilde{z}_{ii} \right] \right]$$

$\nabla f = 0$

$$\Rightarrow \tilde{z}_{ii} - \varepsilon_i + \lambda = 0$$

$$\Rightarrow \tilde{z}_{ii} = \underline{\varepsilon_i - \lambda}$$

$$\tilde{z}_{ii} = [\varepsilon_i - \lambda]_+$$

$$Z = U \cdot \text{diag}([\varepsilon_i - \lambda]_+) U^T$$

Example: graphical lasso

When learning a sparse Gaussian graphical model, one resorts to:

$$\begin{aligned} \text{minimize}_{\Theta} \quad & \underbrace{-\log \det \Theta + \langle \Theta, S \rangle}_{\text{negative log-likelihood of Gaussian graphical model}} + \underbrace{\lambda \|\Theta\|_1}_{\text{encourage sparsity}} \\ \text{s.t.} \quad & \Theta \succeq 0 \end{aligned}$$

$$\begin{aligned} \text{minimize}_{\Theta} \quad & -\log \det \Theta + \langle \Theta, S \rangle + \mathbb{I}_{\mathbb{S}_+}(\Theta) + \lambda \|\Psi\|_1 \\ \text{s.t.} \quad & \Theta = \Psi \end{aligned} \quad (10.3)$$

where $\mathbb{S}_+ := \{X \mid X \succeq 0\}$

$$f(\theta) = \log \det \theta, \quad \theta \geq 0$$

\Downarrow

$$\text{define } g(t) = \log(\det(\theta + tV)), \quad \theta + tV \geq 0$$

$$\begin{aligned} g(t) &= \log \det \left(\theta^{\frac{1}{2}} \cdot \theta^{\frac{1}{2}} + t \cdot \theta^{\frac{1}{2}} \theta^{-\frac{1}{2}} \cdot V \cdot \theta^{-\frac{1}{2}} \cdot \theta^{\frac{1}{2}} \right) \\ &= \log \det \left(\underbrace{\theta^{\frac{1}{2}}}_{\Delta} \left(I + t \underbrace{\theta^{-\frac{1}{2}} V \theta^{-\frac{1}{2}}}_{\Delta} \right) \underbrace{\theta^{\frac{1}{2}}}_{\Delta} \right) \end{aligned}$$

$$\det(AB) = \det(A) \cdot \det(B)$$

$$\begin{aligned} \underline{g(t)} &= \log \det(\theta) + \log \det \left(I + t \underbrace{\theta^{-\frac{1}{2}} V \theta^{-\frac{1}{2}}}_{\lambda_1, \dots, \lambda_n} \right) \\ &= \log \det(\theta) + \sum_{i=1}^n \log(1 + t\lambda_i) \end{aligned}$$

$$g'(t) = \sum_{i=1}^n \frac{\lambda_i}{1 + t\lambda_i}, \quad g''(t) = - \sum_{i=1}^n \frac{\lambda_i^2}{(1 + t\lambda_i)^2} \leq 0$$

$g(t)$ concave, $f(\theta)$ concave

Example: graphical lasso

ADMM for solving (10.3):

$$\Theta^{t+1} = \arg \min_{\Theta \succeq \mathbf{0}} \left\{ -\log \det \Theta + \frac{\rho}{2} \left\| \Theta - \Psi^t + \frac{1}{\rho} \Lambda^t + \frac{1}{\rho} \mathbf{S} \right\|_{\text{F}}^2 \right\}$$

$$\Psi^{t+1} = \arg \min_{\Psi} \left\{ \lambda \|\Psi\|_1 + \frac{\rho}{2} \left\| \Theta^{t+1} - \Psi + \frac{1}{\rho} \Lambda^t \right\|_{\text{F}}^2 \right\}$$

$$\Lambda^{t+1} = \Lambda^t + \rho (\Theta^{t+1} - \Psi^{t+1})$$

Example: graphical lasso

This is equivalent to

$$\begin{aligned}\Theta^{t+1} &= \mathcal{F}_\rho \left(\Psi^t - \frac{1}{\rho} \Lambda^t - \frac{1}{\rho} S \right) \\ \Psi^{t+1} &= \text{ST}_{\lambda\rho^{-1}} \left(\Theta^{t+1} + \frac{1}{\rho} \Lambda^t \right) \quad (\text{soft thresholding}) \\ \Lambda^{t+1} &= \Lambda^t + \rho (\Theta^{t+1} - \Psi^{t+1})\end{aligned}$$

where for $\mathbf{X} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top \succeq \mathbf{0}$ with $\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$, one has

$$\mathcal{F}_\rho(\mathbf{X}) := \frac{1}{2} \mathbf{U} \text{diag}(\{\lambda_i + \sqrt{\lambda_i^2 + \frac{4}{\rho}}\}) \mathbf{U}^\top$$

Example: consensus optimization

Consider solving the following minimization problem

$$\text{minimize}_x \sum_{i=1}^N f_i(x)$$

$$\min_{x \in C_i} f_i(x)$$

$$\text{minimize} \sum_{i=1}^N f_i(x_i) \quad (\text{block separable})$$

$$x_1 \dots x_n$$

$$\text{s.t.} \quad x_i = z \quad 1 \leq i \leq N$$

$$\text{minimize} \sum_{i=1}^N f_i(x_i)$$

$$x_1 = z$$

$$x_2 = z$$

$$\vdots$$

$$x_n = z$$

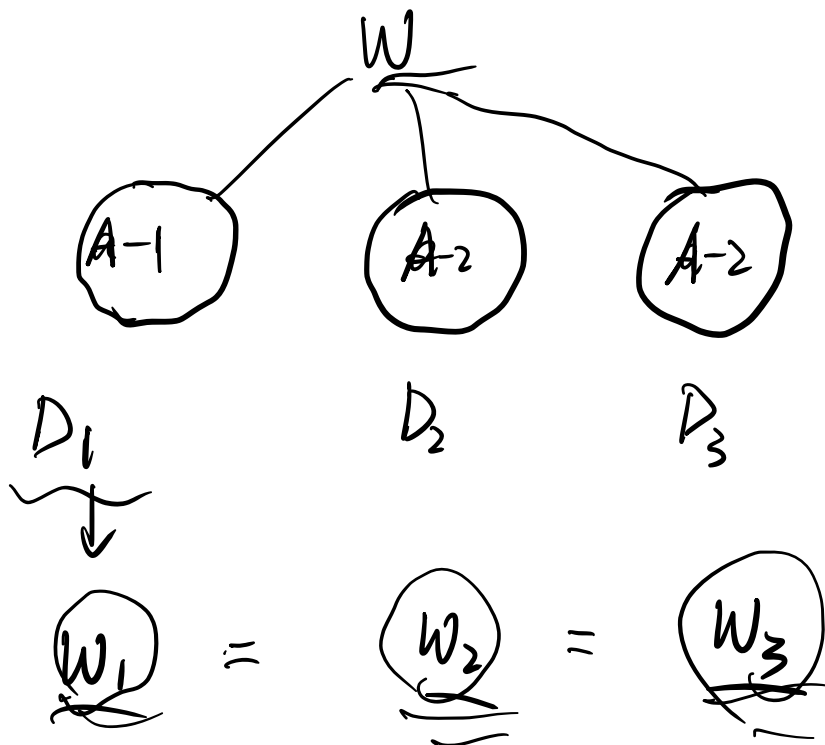
$$\Leftrightarrow \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} I \\ I \\ \vdots \\ I \end{bmatrix} z$$

$\underbrace{\hspace{10em}}_U$

$$Ax + Bz = b$$

$$\bar{I} \cdot u + \begin{bmatrix} I \\ I \\ \vdots \\ I \end{bmatrix} z = 0$$

$$\text{s.t.} \quad \underbrace{u := \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}}_{\text{}} = \underbrace{\begin{bmatrix} I \\ \vdots \\ I \end{bmatrix}}_{\text{}} \underbrace{z}_{\text{}}$$



$$\underline{W_1 = z}$$

Example: consensus optimization

ADMM for solving this problem: \Leftrightarrow

$$\begin{aligned} \mathbf{u}^{t+1} &= \arg \min_{\mathbf{u}=[\mathbf{x}_i]_{1 \leq i \leq N}} \left\{ \underbrace{\sum_{i=1}^N f_i(\mathbf{x}_i)} + \underbrace{\frac{\rho}{2} \sum_{i=1}^N \left\| \mathbf{x}_i - \mathbf{z}^t + \frac{1}{\rho} \boldsymbol{\lambda}_i^t \right\|_2^2} \right\} \checkmark \\ \mathbf{z}^{t+1} &= \arg \min_{\mathbf{z}} \left\{ \underbrace{\frac{\rho}{2} \sum_{i=1}^N \left\| \mathbf{x}_i^{t+1} - \mathbf{z} + \frac{1}{\rho} \boldsymbol{\lambda}_i^t \right\|_2^2}_{2'} \right\} \checkmark \\ \boldsymbol{\lambda}_i^{t+1} &= \boldsymbol{\lambda}_i^t + \rho(\mathbf{x}_i^{t+1} - \mathbf{z}^{t+1}), \quad 1 \leq i \leq N \end{aligned}$$

Example: consensus optimization

This is equivalent to

$$\mathbf{x}_i^{t+1} = \arg \min_{\mathbf{x}_i} \left\{ f_i(\mathbf{x}_i) + \frac{\rho}{2} \left\| \mathbf{x}_i - \mathbf{z}^t + \frac{1}{\rho} \boldsymbol{\lambda}_i^t \right\|_2^2 \right\} \quad 1 \leq i \leq N$$

(can be computed in parallel)

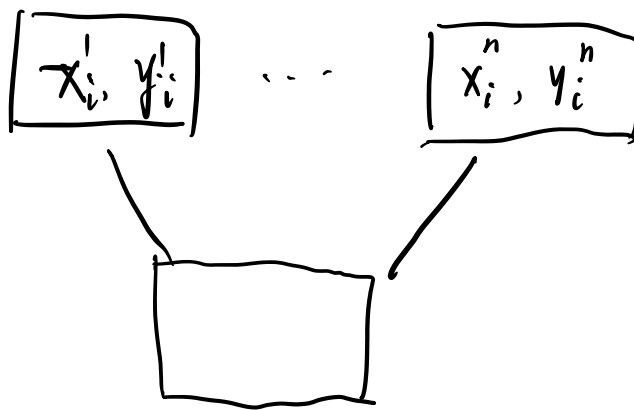
$$\mathbf{z}^{t+1} = \frac{1}{N} \sum_{i=1}^N \left(\mathbf{x}_i^{t+1} + \frac{1}{\rho} \boldsymbol{\lambda}_i^t \right)$$

(gather all local iterates)

$$\boldsymbol{\lambda}_i^{t+1} = \boldsymbol{\lambda}_i^t + \rho(\mathbf{x}_i^{t+1} - \mathbf{z}^{t+1}), \quad 1 \leq i \leq N$$

(“broadcast” \mathbf{z}^{t+1} to update all local multipliers)

ADMM is well suited for distributed optimization!



$$\arg \min_{W^1} \| W^1 x^1 - y^1 \|_2^2$$

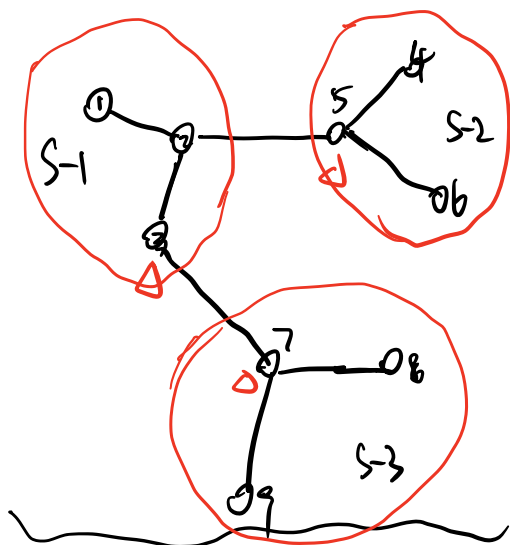
$$\arg \min_{W^n} \| W^n x^n - y^n \|_2^2$$

~~~~~

$$\sum_{k=1}^n \| W^k x^k - y^k \|_2^2$$

$$\text{s.t. } W^k = z, \quad k=1, \dots, n$$

~~~~~



Centralized
1. Server to control

$$\min_{x_1, \dots, x_9} f(x_1, \dots, x_9)$$

$$\min f_1(x_1, x_2, x_3) + f_2(x_4, x_5, x_6) + f_3(x_7, x_8, x_9)$$

$$\underline{S'_1: 1. 2. 3. 5. 7}$$

$$\underline{S'_2: 4. 5. 6. 2}$$

$$\underline{S'_3: 7. 8. 9. 3}$$

$$\min_{\Delta \quad \Delta} f_1(x_1, x_2, x_3, x_5, x_7) -$$

$$\min_{\Delta} f_2(x_4, x_5, x_6, x_2) -$$

$$\min_{\Delta} f_3(x_7, x_8, x_9, x_3) -$$

$$\min f_1(x_1, x_2, x_3) \quad \text{---} \quad \checkmark$$

$$\min f_2(x_4, x_5, x_6) \quad \text{---} \quad \checkmark$$

$$\min f_3(x_7, x_8, x_9) \quad \text{---} \quad \checkmark$$

$$\underline{\textcircled{Z}: 1-9}$$

Convergence of ADMM

Theorem 10.1 (Convergence of ADMM)

Suppose f_1 and f_2 are closed convex functions, and γ is any constant obeying $\gamma \geq 2\|\boldsymbol{\lambda}^*\|_2$. Then

$$F(\mathbf{x}^{(t)}, \mathbf{z}^{(t)}) - F^{\text{opt}} \leq \frac{\|\mathbf{z}^0 - \mathbf{z}^*\|_{\rho \mathbf{B}^\top \mathbf{B}}^2 + \frac{(\gamma + \|\boldsymbol{\lambda}^0\|_2)^2}{\rho}}{2(t+1)} \quad (10.4a)$$

$$\|\mathbf{A}\mathbf{x}^{(t)} + \mathbf{B}\mathbf{z}^{(t)} - \mathbf{b}\|_2 \leq \frac{\|\mathbf{z}^0 - \mathbf{z}^*\|_{\rho \mathbf{B}^\top \mathbf{B}}^2 + \frac{(\gamma + \|\boldsymbol{\lambda}^0\|_2)^2}{\rho}}{\gamma(t+1)} \quad (10.4b)$$

where $\mathbf{x}^{(t)} := \frac{1}{t+1} \sum_{k=1}^{t+1} \mathbf{x}^k$, $\mathbf{z}^{(t)} := \frac{1}{t+1} \sum_{k=1}^{t+1} \mathbf{z}^k$, and for any \mathbf{C} , $\|\mathbf{z}\|_{\mathbf{C}}^2 := \mathbf{z}^\top \mathbf{C} \mathbf{z}$

- convergence rate: $O(1/t)$
- iteration complexity: $O(1/\varepsilon)$

Fundamental inequality

Define

$$\mathbf{w} := \begin{bmatrix} x \\ z \\ \lambda \end{bmatrix}, \quad \mathbf{w}^t := \begin{bmatrix} x^t \\ z^t \\ \lambda^t \end{bmatrix}, \quad \mathbf{G} := \begin{bmatrix} & \mathbf{A}^\top & \\ -\mathbf{A} & -\mathbf{B} & \mathbf{B}^\top \end{bmatrix}, \quad \mathbf{d} := \begin{bmatrix} 0 \\ 0 \\ b \end{bmatrix}$$
$$\mathbf{H} := \begin{bmatrix} 0 & & \\ & \rho \mathbf{B}^\top \mathbf{B} & \\ & & \rho^{-1} \mathbf{I} \end{bmatrix}, \quad \|\mathbf{w}\|_H^2 := \mathbf{w}^\top \mathbf{H} \mathbf{w}$$

Lemma 10.2

For any x, z, λ , one has

$$\begin{aligned} F(x, z) - F(x^{t+1}, z^{t+1}) + \langle \mathbf{w} - \mathbf{w}^{t+1}, \mathbf{G}\mathbf{w} + \mathbf{d} \rangle \\ \geq \underbrace{\frac{1}{2} \|\mathbf{w} - \mathbf{w}^{t+1}\|_H^2 - \frac{1}{2} \|\mathbf{w} - \mathbf{w}^t\|_H^2} \end{aligned}$$

$$Gw+d = \begin{bmatrix} 0 & 0 & A^T \\ 0 & 0 & B^T \\ A-B & 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ z \\ \lambda \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ b \end{bmatrix}$$

$$= \begin{bmatrix} A^T \lambda \\ B^T \lambda \\ -Ax - Bz + b \end{bmatrix}$$

$$\langle \overset{t}{w-w}, Gw+d \rangle = \left\langle \begin{bmatrix} x^{(t)} - x \\ z^{(t)} - z \\ \lambda^{(t)} - \lambda \end{bmatrix}, \begin{bmatrix} A^T \lambda \\ B^T \lambda \\ -Ax - Bz + b \end{bmatrix} \right\rangle$$

$$= (x^{(t)T} - x^T) (A^T \lambda) + (z^{(t)T} - z^T) (B^T \lambda)$$

$$+ (\lambda^{(t)T} - \lambda^T) (-Ax - Bz + b)$$

$$= \lambda^{(t)T} (-Ax - Bz + b) + \langle \lambda, Ax + Bz - b \rangle$$

Proof of Theorem 10.1

Set $\mathbf{x} = \mathbf{x}^*$, $\mathbf{z} = \mathbf{z}^*$, and $\mathbf{w} = [\mathbf{x}^{*\top}, \mathbf{z}^{*\top}, \boldsymbol{\lambda}^\top]^\top$ in Lemma 10.2 to reach

$$F(\mathbf{x}^*, \mathbf{z}^*) - F(\mathbf{x}^{k+1}, \mathbf{z}^{k+1}) + \langle \mathbf{w} - \mathbf{w}^{k+1}, \mathbf{G}\mathbf{w} + \mathbf{d} \rangle \geq \underbrace{\frac{\|\mathbf{w} - \mathbf{w}^{k+1}\|_H^2}{2} - \frac{\|\mathbf{w} - \mathbf{w}^k\|_H^2}{2}}_{\text{forms telescopic sum}}$$

Summing over all $k = 0, \dots, t$ gives

$$\begin{aligned} & (t+1)F(\mathbf{x}^*, \mathbf{z}^*) - \sum_{k=1}^{t+1} F(\mathbf{x}^k, \mathbf{z}^k) + \left\langle (t+1)\mathbf{w} - \sum_{k=1}^{t+1} \mathbf{w}^k, \mathbf{G}\mathbf{w} + \mathbf{d} \right\rangle \\ & \geq \frac{\|\mathbf{w} - \mathbf{w}^{t+1}\|_H^2 - \|\mathbf{w} - \mathbf{w}^0\|_H^2}{2} \end{aligned}$$

If we define

$$\mathbf{w}^{(t)} = \frac{1}{t+1} \sum_{k=1}^{t+1} \mathbf{w}^k, \quad \mathbf{x}^{(t)} = \frac{1}{t+1} \sum_{k=1}^{t+1} \mathbf{x}^k, \quad \mathbf{z}^{(t)} = \frac{1}{t+1} \sum_{k=1}^{t+1} \mathbf{z}^k, \quad \boldsymbol{\lambda}^{(t)} = \frac{1}{t+1} \sum_{k=1}^{t+1} \boldsymbol{\lambda}^k,$$

then from convexity of F we have

$$\underbrace{F(\mathbf{x}^{(t)}, \mathbf{z}^{(t)})}_{\text{ADMM}} - \underbrace{F(\mathbf{x}^*, \mathbf{z}^*)}_{= F^{\text{opt}}} + \underbrace{\langle \mathbf{w}^{(t)} - \mathbf{w}, \mathbf{G}\mathbf{w} + \mathbf{d} \rangle}_{\text{ADMM}} \leq \frac{1}{2(t+1)} \|\mathbf{w} - \mathbf{w}^0\|_H^2$$

$$F\left(\frac{1}{t+1} \sum_{k=1}^{t+1} x^k, \frac{1}{t+1} \sum_{k=1}^{t+1} z^k\right)$$

$$= F(x^{(t)}, z^{(t)})$$

$$\leq \frac{1}{t+1} \sum_{k=1}^{t+1} F(x^k, z^k)$$

~~~~~

# Proof of Theorem 10.1

Further, we claim that

$$\langle \mathbf{w}^{(t)} - \mathbf{w}, \mathbf{G}\mathbf{w} + \mathbf{d} \rangle = \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x}^{(t)} + \mathbf{B}\mathbf{z}^{(t)} - \mathbf{b} \rangle \quad (10.5)$$

which together with preceding bounds yields

$$\begin{aligned} F(\mathbf{x}^{(t)}, \mathbf{z}^{(t)}) - F^{\text{opt}} + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x}^{(t)} + \mathbf{B}\mathbf{z}^{(t)} - \mathbf{b} \rangle &\leq \frac{1}{2(t+1)} \|\mathbf{w} - \mathbf{w}^0\|_H^2 \\ &= \frac{1}{2(t+1)} \left\{ \|\mathbf{z} - \mathbf{z}^0\|_{\rho \mathbf{B}^\top \mathbf{B}}^2 + \frac{1}{\rho} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^0\|_2^2 \right\} \end{aligned}$$

Notably, this holds for any  $\boldsymbol{\lambda}$

Taking maximum of both sides over  $\{\boldsymbol{\lambda} \mid \|\boldsymbol{\lambda}\|_2 \leq \gamma\}$  yields

$$\begin{aligned} F(\mathbf{x}^{(t)}, \mathbf{z}^{(t)}) - F^{\text{opt}} + \gamma \|\mathbf{A}\mathbf{x}^{(t)} + \mathbf{B}\mathbf{z}^{(t)} - \mathbf{b}\|_2 &\leq \underbrace{(\gamma + \|\boldsymbol{\lambda}^0\|_2)^2}_{\leq (\gamma + \|\boldsymbol{\lambda}^0\|_2)^2} \\ &\leq \frac{\left\{ \|\mathbf{z} - \mathbf{z}^0\|_{\rho \mathbf{B}^\top \mathbf{B}}^2 + \frac{(\gamma + \|\boldsymbol{\lambda}^0\|_2)^2}{\rho} \right\}}{2(t+1)} \end{aligned} \quad (10.6)$$

which immediately establishes (10.4a)

## Proof of Theorem 10.1 (cont.)

---

Caution needs to be exercised since, in general, (10.6) does not establish (10.4b), since  $F(\mathbf{x}^{(t)}, \mathbf{z}^{(t)}) - F^{\text{opt}}$  may be negative (as  $(\mathbf{x}^{(t)}, \mathbf{z}^{(t)})$  is not guaranteed to be feasible)

Fortunately, if  $\gamma \geq 2\|\boldsymbol{\lambda}^*\|_2$ , then standard results (e.g. Theorem 3.60 in Beck '18) reveal that  $F(\mathbf{x}^{(t)}, \mathbf{z}^{(t)}) - F^{\text{opt}}$  will not be “too negative”, thus completing proof

# Proof of Theorem 10.1

---

Finally, we prove (10.5). Observe that

$$\begin{aligned}
 \langle w^{(t)} - w, Gw + d \rangle &= \underbrace{\langle w^{(t)} - w, G(w - w^{(t)}) \rangle}_{=0 \text{ since } G \text{ is skew-symmetric}} + \langle w^{(t)} - w, Gw^{(t)} + d \rangle \\
 &= \langle w^{(t)} - w, Gw^{(t)} + d \rangle
 \end{aligned} \tag{10.7}$$

To further simplify this inner product, we use  $Ax^* + Bz^* = b$  to obtain

$$\begin{aligned}
 \langle w^{(t)} - w, Gw^{(t)} + d \rangle &= \langle x^{(t)} - x^*, A^\top \lambda^{(t)} \rangle + \langle z^{(t)} - z^*, B^\top \lambda^{(t)} \rangle \\
 &\quad + \langle \lambda^{(t)} - \lambda, -Ax^{(t)} - Bz^{(t)} + b \rangle \\
 &= \langle -Ax^* - Bz^* + b, \lambda^{(t)} \rangle + \langle \lambda, Ax^{(t)} + Bz^{(t)} - b \rangle \\
 &= \langle \lambda, Ax^{(t)} + Bz^{(t)} - b \rangle
 \end{aligned}$$

# Proof of Lemma 10.2

---

To begin with, ADMM update rule requires

$$-\rho \mathbf{A}^\top \left( \mathbf{A} \mathbf{x}^{t+1} + \underbrace{\mathbf{B} \mathbf{z}^t}_{\text{ADMM}} - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right) \in \partial f_1(\mathbf{x}^{t+1}) \quad \checkmark$$

$$-\rho \mathbf{B}^\top \left( \mathbf{A} \mathbf{x}^{t+1} + \mathbf{B} \mathbf{z}^{t+1} - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right) \in \partial f_2(\mathbf{z}^{t+1}) \quad \checkmark$$

Therefore, for any  $\mathbf{x}, \mathbf{z}$ ,

$$f_1(\mathbf{x}) - f_1(\mathbf{x}^{t+1}) + \left\langle \rho \mathbf{A}^\top \left( \mathbf{A} \mathbf{x}^{t+1} + \underbrace{\mathbf{B} \mathbf{z}^t}_{\text{ADMM}} - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right), \mathbf{x} - \mathbf{x}^{t+1} \right\rangle \geq 0 \quad \checkmark_{\Delta}$$

$$f_2(\mathbf{z}) - f_2(\mathbf{z}^{t+1}) + \left\langle \rho \mathbf{B}^\top \left( \mathbf{A} \mathbf{x}^{t+1} + \underbrace{\mathbf{B} \mathbf{z}^{t+1}}_{\text{ADMM}} - \mathbf{b} + \frac{1}{\rho} \boldsymbol{\lambda}^t \right), \mathbf{z} - \mathbf{z}^{t+1} \right\rangle \geq 0 \quad \checkmark_{\Delta}$$



# Proof of Lemma 10.2 (cont.)

Using  $\underline{\lambda}^{t+1} = \lambda^t + \rho(\underline{Ax}^{t+1} + \underline{Bz}^{t+1} - b)$ , setting  $\underline{\tilde{\lambda}}^t := \lambda^t + \rho(\underline{Ax}^{t+1} + \underline{Bz}^t - b)$ , and adding above two inequalities give

$$F(x, z) - F(x^{t+1}, z^{t+1}) + \left\langle \begin{bmatrix} x - x^{t+1} \\ z - z^{t+1} \\ \lambda - \tilde{\lambda}^t \end{bmatrix}, \begin{bmatrix} A^\top \tilde{\lambda}^t \\ B^\top \tilde{\lambda}^t \\ -Ax^{t+1} - Bz^{t+1} + b \end{bmatrix} - \begin{bmatrix} 0 \\ \rho B^\top B(z^t - z^{t+1}) \\ \frac{1}{\rho}(\lambda^t - \lambda^{t+1}) \end{bmatrix} \right\rangle \geq 0 \quad (10.8)$$

Next, we'd like to simplify above inner product. Let  $C := \rho B^\top B$ , then

$$(z - z^{t+1})^\top C (z^t - z^{t+1}) = \frac{1}{2} \|z - z^{t+1}\|_C^2 - \frac{1}{2} \|z - z^t\|_C^2 + \frac{1}{2} \|z^t - z^{t+1}\|_C^2$$

$$(a-b)^\top (c-f) =$$

# Proof of Lemma 10.2 (cont.)

---

Also,

$$\begin{aligned}
 & 2(\boldsymbol{\lambda} - \boldsymbol{\lambda}^{t+1})^\top (\boldsymbol{\lambda}^t - \boldsymbol{\lambda}^{t+1}) = \\
 & = \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^{t+1}\|_2^2 - \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^t\|_2^2 + \|\tilde{\boldsymbol{\lambda}}^t - \boldsymbol{\lambda}^t\|_2^2 - \|\tilde{\boldsymbol{\lambda}}^t - \boldsymbol{\lambda}^{t+1}\|_2^2 \\
 & = \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^{t+1}\|_2^2 - \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^t\|_2^2 + \rho^2 \|\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^t - \mathbf{b}\|_2^2 \\
 & \quad - \|\boldsymbol{\lambda}^t + \rho(\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^t - \mathbf{b}) - \boldsymbol{\lambda}^t - \rho(\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^{t+1} - \mathbf{b})\|_2^2 \\
 & = \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^{t+1}\|_2^2 - \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^t\|_2^2 + \rho^2 \underbrace{\|\mathbf{A}\mathbf{x}^{t+1} + \mathbf{B}\mathbf{z}^t - \mathbf{b}\|_2^2}_{\text{}} \\
 & \quad - \rho^2 \|\mathbf{B}(\mathbf{z}^t - \mathbf{z}^{t+1})\|_2^2 \quad \checkmark
 \end{aligned}$$

which implies that

$$\begin{aligned}
 & 2(\boldsymbol{\lambda} - \boldsymbol{\lambda}^{t+1})^\top (\boldsymbol{\lambda}^t - \boldsymbol{\lambda}^{t+1}) \\
 & \geq \underbrace{\|\boldsymbol{\lambda} - \boldsymbol{\lambda}^{t+1}\|_2^2 - \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^t\|_2^2}_{\text{}} - \rho^2 \|\mathbf{B}(\mathbf{z}^t - \mathbf{z}^{t+1})\|_2^2
 \end{aligned}$$

## Proof of Lemma 10.2 (cont.)

---

Combining above results gives

$$\begin{aligned}
 & \left\langle \begin{bmatrix} \mathbf{x} - \mathbf{x}^{t+1} \\ \mathbf{z} - \mathbf{z}^{t+1} \\ \lambda - \tilde{\lambda}^t \end{bmatrix}, \begin{bmatrix} \mathbf{0} \\ \rho \mathbf{B}^\top \mathbf{B}(\mathbf{z}^t - \mathbf{z}^{t+1}) \\ \frac{1}{\rho}(\lambda^t - \lambda^{t+1}) \end{bmatrix} \right\rangle \\
 & \geq \frac{1}{2} \|\mathbf{w} - \mathbf{w}^{t+1}\|_H^2 - \frac{1}{2} \|\mathbf{w} - \mathbf{w}^t\|_H^2 + \frac{1}{2} \|\mathbf{z}^t - \mathbf{z}^{t+1}\|_C^2 - \frac{\rho}{2} \|\mathbf{B}(\mathbf{z}^t - \mathbf{z}^{t+1})\|_2^2 \quad \checkmark \\
 & = \frac{1}{2} \|\mathbf{w} - \mathbf{w}^{t+1}\|_H^2 - \frac{1}{2} \|\mathbf{w} - \mathbf{w}^t\|_H^2 \quad \checkmark
 \end{aligned}$$

This together with (10.8) yields

$$\begin{aligned}
 & F(\mathbf{x}, \mathbf{z}) - F(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}) + \langle \mathbf{w} - \mathbf{w}^{t+1}, \mathbf{G}\mathbf{w}^{t+1} + \mathbf{d} \rangle \\
 & \geq \frac{1}{2} \|\mathbf{w} - \mathbf{w}^{t+1}\|_H^2 - \frac{1}{2} \|\mathbf{w} - \mathbf{w}^t\|_H^2 \quad \checkmark
 \end{aligned}$$

Since  $\mathbf{G}$  is skew-symmetric, repeating prior argument in (10.7) gives

$$\langle \mathbf{w} - \mathbf{w}^{t+1}, \mathbf{G}\mathbf{w}^{t+1} + \mathbf{d} \rangle = \langle \mathbf{w} - \mathbf{w}^{t+1}, \mathbf{G}\mathbf{w} + \mathbf{d} \rangle \quad \checkmark$$

This immediately completes proof

# Convergence of ADMM in practice

---

$$\underbrace{O\left(\frac{1}{\epsilon}\right)}$$

- ADMM is slow to converge to high accuracy
- ADMM often converges to modest accuracy within a few tens of iterations, which is sufficient for many large-scale applications

# Beyond two-block models

---

Convergence is not guaranteed when there are 3 or more blocks

- e.g. consider solving

$$x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + x_3 \mathbf{a}_3 = \mathbf{0}$$

where

$$[\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3] = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 2 & 2 \end{bmatrix}$$

3-block ADMM is divergent for solving this problem (Chen et al. '16)



# Reference

---

- [1] "*Distributed optimization and statistical learning via the alternating direction method of multipliers*," S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *Foundations and Trends in Machine learning*, 2011.
- [2] "*A First Course in Convex Optimization Theory*," E. Ryu, W. Yin.
- [3] "*First-order methods in optimization*," A. Beck, Vol. 25, SIAM, 2017.
- [4] "*Mathematical optimization, MATH301 lecture notes*," E. Candes, Stanford.
- [5] "*Optimization methods for large-scale systems, EE236C lecture notes*," L. Vandenberghe, UCLA.
- [6] "*Large scale optimization for machine learning, ISE633 lecture notes*," M. Razaviyayn, USC.
- [7] "*Modern big data optimization, IE487/587 lecture notes*," M. Hong, ISU.

# Reference

---

- [8] "*Convex optimization and algorithms*," D. Bertsekas, 2015.
- [9] "*Proximal algorithms*," N. Parikh and S. Boyd, *Foundations and Trends in Optimization*, 2013.
- [10] "*The direct extension of ADMM for multi-block convex minimization problems is not necessarily convergent*," C. Chen, B. He, Y. Ye, X. Yuan, *Mathematical Programming*, 2016.
- [11] "*Robust principal component analysis?*," E. Candes, X. Li, Y. Ma, J. Wright, *Journal of the ACM*, 2011.
- [12] "*Sparse inverse covariance estimation with the graphical lasso*," J. Friedman, T. Hastie, and R. Tibshirani, *Biostatistics*, 2008.