#### **Distributed Optimization**

Optimization Techniques (ENGG\*6140)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Winter 2023 **Problem Statement** 

• When we have several optimization variables, we have:

$$\min_{\substack{\{\mathbf{x}_i\}_{i=1}^m}} f(\mathbf{x}_1, \dots, \mathbf{x}_m), \tag{1}$$

where m is the number of optimization variables.

In this case, we can have distributed optimization because we can work on the
optimization variables in a distributed manner.

Alternating Optimization

## Alternating Optimization

• Consider the following multivariate optimization problem:

$$\min_{\{\boldsymbol{x}_i\}_{i=1}^m} \quad f(\boldsymbol{x}_1,\ldots,\boldsymbol{x}_m),$$

where the objective function depends on m variables.

- When we have several optimization variables, we can alternate between optimizing over each of these variables. This technique is called **alternating optimization** in the literature [1] (also see [2, Chapter 4]).
- Alternating optimization alternates between updating every variable while assuming other variables are constant, set to their last updated value. After random feasible initialization, it updates solutions as [1]:

$$\begin{split} \mathbf{x}_{1}^{(k+1)} &:= \arg\min_{\mathbf{x}_{1}} f(\mathbf{x}_{1}, \mathbf{x}_{2}^{(k)}, \dots, \mathbf{x}_{m-1}^{(k)}, \mathbf{x}_{m}^{(k)}), \\ \mathbf{x}_{2}^{(k+1)} &:= \arg\min_{\mathbf{x}_{2}} f(\mathbf{x}_{1}^{(k+1)}, \mathbf{x}_{2}, \dots, \mathbf{x}_{m-1}^{(k)}, \mathbf{x}_{m}^{(k)}), \\ \vdots \\ \mathbf{x}_{m}^{(k+1)} &:= \arg\min_{\mathbf{x}_{m}} f(\mathbf{x}_{1}^{(k+1)}, \mathbf{x}_{2}^{(k+1)}, \dots, \mathbf{x}_{m-1}^{(k+1)}, \mathbf{x}_{m}), \end{split}$$

until convergence.

- Any optimization methods, including first-order and second-order methods, can be used for each of the optimization lines above.
- In most cases, alternating optimization is robust to changing the order of updates of variables.

# Alternating Optimization for Decomposable Function in terms of Variables

• If the function  $f(x_1, \ldots, x_m)$  is decomposable in terms of variables, i.e., if we have:

$$f(\mathbf{x}_1,\ldots,\mathbf{x}_m)=\sum_{i=1}^m f_i(\mathbf{x}_i),$$

the alternating optimization can be simplified to:

$$\begin{aligned} \mathbf{x}_{1}^{(k+1)} &:= \arg\min_{\mathbf{x}_{1}} f_{1}(\mathbf{x}_{1}), \\ \mathbf{x}_{2}^{(k+1)} &:= \arg\min_{\mathbf{x}_{2}} f_{2}(\mathbf{x}_{2}), \\ &\vdots \\ &\mathbf{x}_{m}^{(k+1)} &:= \arg\min_{\mathbf{x}_{m}} f_{m}(\mathbf{x}_{m}), \end{aligned}$$

because other terms become constant in optimization.

- The above updates mean that if the function is completely decomposable in terms of variables, the updates of variables are independent and can be done independently.
- Hence, in that case, alternating optimization is reduced to *m* independent optimization problems, each of which can be solved by any optimization method such as the first-order and second-order methods.

#### Proximal Alternating Optimization

• Proximal alternating optimization uses proximal operator:

$$\operatorname{prox}_{\lambda g}(\mathbf{x}) := \arg\min_{\mathbf{u}} \left( g(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{x}\|_2^2 \right), \tag{2}$$

for minimization to keep the updated solution close to the solution of previous iteration [1]:

$$\begin{aligned} \mathbf{x}_{1}^{(k+1)} &:= \arg\min_{\mathbf{x}_{1}} \left( f(\mathbf{x}_{1}, \mathbf{x}_{2}^{(k)}, \dots, \mathbf{x}_{m-1}^{(k)}, \mathbf{x}_{m}^{(k)}) + \frac{1}{2\lambda} \|\mathbf{x}_{1} - \mathbf{x}_{1}^{(k)}\|_{2}^{2} \right), \\ \mathbf{x}_{2}^{(k+1)} &:= \arg\min_{\mathbf{x}_{2}} \left( f(\mathbf{x}_{1}^{(k+1)}, \mathbf{x}_{2}, \dots, \mathbf{x}_{m-1}^{(k)}, \mathbf{x}_{m}^{(k)}) + \frac{1}{2\lambda} \|\mathbf{x}_{2} - \mathbf{x}_{2}^{(k)}\|_{2}^{2} \right), \\ \vdots \\ \mathbf{x}_{m}^{(k+1)} &:= \arg\min_{\mathbf{x}_{m}} \left( f(\mathbf{x}_{1}^{(k+1)}, \mathbf{x}_{2}^{(k+1)}, \dots, \mathbf{x}_{m-1}^{(k+1)}, \mathbf{x}_{m}) + \frac{1}{2\lambda} \|\mathbf{x}_{m} - \mathbf{x}_{m}^{(k)}\|_{2}^{2} \right). \end{aligned}$$

# Alternating Optimization for Constrained Problems

• The alternating optimization methods can also be used for constrained problems:

$$\begin{array}{l} \underset{\{\boldsymbol{x}_i\}_{i=1}^m}{\min } \quad f(\boldsymbol{x}_1, \dots, \boldsymbol{x}_m) \\ \text{subject to} \quad \boldsymbol{x}_i \in \mathcal{S}_i, \quad \forall i \in \{1, \dots, m\}. \end{array}$$
(3)

• In this case, every line of the optimization is a constrained problem:

$$\begin{aligned} \mathbf{x}_{1}^{(k+1)} &:= \arg\min_{\mathbf{x}_{1}} \left( f(\mathbf{x}_{1}, \mathbf{x}_{2}^{(k)}, \dots, \mathbf{x}_{m-1}^{(k)}, \mathbf{x}_{m}^{(k)}), \text{ s.t. } \mathbf{x}_{1} \in \mathcal{S}_{1} \right), \\ \mathbf{x}_{2}^{(k+1)} &:= \arg\min_{\mathbf{x}_{2}} \left( f(\mathbf{x}_{1}^{(k+1)}, \mathbf{x}_{2}, \dots, \mathbf{x}_{m-1}^{(k)}, \mathbf{x}_{m}^{(k)}), \text{ s.t. } \mathbf{x}_{2} \in \mathcal{S}_{2} \right), \\ &\vdots \\ \mathbf{x}_{m}^{(k+1)} &:= \arg\min_{\mathbf{x}_{m}} \left( f(\mathbf{x}_{1}^{(k+1)}, \mathbf{x}_{2}^{(k+1)}, \dots, \mathbf{x}_{m-1}^{(k+1)}, \mathbf{x}_{m}), \text{ s.t. } \mathbf{x}_{m} \in \mathcal{S}_{m} \right) \end{aligned}$$

- Any constrained optimization methods can be used for each of the optimization lines above. Some examples are projected gradient method, proximal methods, interior-point methods, etc.
- Practical experiments have shown there is usually no need to use a complete optimization until convergence for every step in the alternating optimization, either unconstrained or constrained. Often, a single step of updating, such as a step of gradient descent or projected gradient method, is enough for the whole algorithm to work.

Dual Ascent and Dual Decomposition Methods

#### **Dual Ascent Method**

Consider the following problem:

$$\begin{array}{ll} \underset{x}{\text{minimize}} & f(x) \\ \text{subject to} & \mathbf{A} \mathbf{x} = \mathbf{b}. \end{array} \tag{4}$$

- We follow the method of Lagrange multipliers discussed before.
- The Lagrangian is:

$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{\nu}) = f(\boldsymbol{x}) + \boldsymbol{\nu}^{\top} (\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}).$$

The dual function is:

$$g(\boldsymbol{\nu}) = \inf_{\boldsymbol{x}} \mathcal{L}(\boldsymbol{x}, \boldsymbol{\nu}). \tag{5}$$

• The optimal dual problem maximizes  $g(\nu)$ :

$$\boldsymbol{\nu}^* = \arg\max_{\boldsymbol{\nu}} g(\boldsymbol{\nu}), \tag{6}$$

so the optimal primal variable is:

$$\boldsymbol{x}^* = \arg\min_{\boldsymbol{x}} \mathcal{L}(\boldsymbol{x}, \boldsymbol{\nu}^*). \tag{7}$$

• For solving Eq. (6), we should take the derivative of the dual function w.r.t. the dual variable:

$$\nabla_{\boldsymbol{\nu}} g(\boldsymbol{\nu}) \stackrel{(5)}{=} \nabla_{\boldsymbol{\nu}} (\inf_{\boldsymbol{x}} \mathcal{L}(\boldsymbol{x}, \boldsymbol{\nu})) \stackrel{(7)}{=} \nabla_{\boldsymbol{\nu}} (f(\boldsymbol{x}^*) + \boldsymbol{\nu}^\top (\boldsymbol{A} \boldsymbol{x}^* - \boldsymbol{b})) = \boldsymbol{A} \boldsymbol{x}^* - \boldsymbol{b}.$$

# Dual Ascent Method

We found:

$$egin{aligned} & m{
u}^* = rg\max_{m{
u}} g(m{
u}), \ & m{x}^* = rg\min_{m{x}} \mathcal{L}(m{x},m{
u}^*). \end{aligned}$$

• The dual problem is a maximization problem so we can use gradient ascent for iteratively updating the dual variable with this gradient. We can alternate between updating the optimal primal and dual variables:

$$\mathbf{x}^{(k+1)} := \arg\min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}^{(k)}), \tag{8}$$

$$\boldsymbol{\nu}^{(k+1)} := \boldsymbol{\nu}^{(k)} + \eta^{(k)} (\boldsymbol{A} \boldsymbol{x}^{(k+1)} - \boldsymbol{b}), \tag{9}$$

where k is the iteration index and  $\eta^{(k)}$  is the step size (also called the learning rate) at iteration k.

- Eq. (8) can be performed by any optimization method. We compute the gradient of *L*(x, ν<sup>(k)</sup>) w.r.t. x. If setting this gradient to zero does not give x in closed form, we can use gradient descent to perform Eq. (8).
- Some papers approximate Eq. (8) by one step or few steps of gradient descent rather than a complete gradient descent until convergence. If using one step, we can write Eq. (8) as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \gamma \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}^{(k)}), \tag{10}$$

where  $\gamma > 0$  is the step size. It has been shown empirically that even one step of gradient descent for Eq. (8) works properly for the whole alternating algorithm.

#### **Dual Ascent Method**

We had:

$$\begin{split} \mathbf{x}^{(k+1)} &:= \arg\min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}^{(k)}), \\ \mathbf{\nu}^{(k+1)} &:= \boldsymbol{\nu}^{(k)} + \eta^{(k)} (\mathbf{A} \mathbf{x}^{(k+1)} - \mathbf{b}), \end{split}$$

- We continue the iterations until convergence of the primal and dual variables to stable values. When we get closer to convergence, we will have  $(\mathbf{A}\mathbf{x}^{k+1} \mathbf{b}) \rightarrow 0$  so that we will not have update of dual variable according to Eq. (9). This means that after convergence, we have  $(\mathbf{A}\mathbf{x}^{k+1} \mathbf{b}) \approx 0$  so that the constraint  $\mathbf{A}\mathbf{x} = \mathbf{b}$  in Eq. (4) is getting satisfied. In other words, the update of dual variable in Eq. (9) is taking care of satisfying the constraint.
- This method is known as the **dual ascent** method because it uses gradient ascent for updating the dual variable.

#### Dual Decomposition Method

• Again, consider the following problem:

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

• If the objective function can be distributed and decomposed on b blocks  $\{x_i\}_{i=1}^{b}$ , i.e.:

$$f(\mathbf{x}) = f_1(\mathbf{x}_1) + \cdots + f_1(\mathbf{x}_b),$$

we can have b Lagrangian functions where the total Lagrangian is the summation of these functions:

$$\mathcal{L}_i(\mathbf{x}_i, \mathbf{\nu}) = f(\mathbf{x}_i) + \mathbf{\nu}^\top (\mathbf{A}\mathbf{x}_i - \mathbf{b}),$$
  
 $\mathcal{L}(\mathbf{x}_i, \mathbf{\nu}) = \sum_{i=1}^{b} (f(\mathbf{x}_i) + \mathbf{\nu}^\top (\mathbf{A}\mathbf{x}_i - \mathbf{b})).$ 

We can divide the Eq. (8), x<sup>(k+1)</sup> := arg min<sub>x</sub> L(x, ν<sup>(k)</sup>), into b updates, each for one of the blocks.

$$\mathbf{x}_{i}^{(k+1)} := \arg\min_{\mathbf{x}_{i}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}^{(k)}), \quad \forall i \in \{1, \dots, b\},$$

$$(11)$$

$$\boldsymbol{\nu}^{(k+1)} := \boldsymbol{\nu}^{(k)} + \eta^{(k)} (\boldsymbol{A} \boldsymbol{x}^{(k+1)} - \boldsymbol{b}).$$
(12)

# Dual Decomposition Method

We found:

$$\begin{split} \mathbf{x}_{i}^{(k+1)} &:= \arg\min_{\mathbf{x}_{i}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}^{(k)}), \quad \forall i \in \{1, \dots, b\}, \\ \mathbf{\nu}^{(k+1)} &:= \boldsymbol{\nu}^{(k)} + \eta^{(k)} (\mathbf{A} \mathbf{x}^{(k+1)} - \mathbf{b}). \end{split}$$

- This is called **dual decomposition** developed by decomposition techniques such as the **Dantzig-Wolfe decomposition** (1960) [3], **Bender's decomposition** (1962) [4], and **Lagrangian decomposition** (1963) [5].
- The dual decomposition methods can divide a problem into sub-problems and solve them in parallel. Hence, it can be used for big data but they are usually slow to converge.

Augmented Lagrangian Method (Method of Multipliers)

# Augmented Lagrangian Method (Method of Multipliers)

• Recall Eq. (4):

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

Assume we regularize the objective function in Eq. (4) by a penalty on not satisfying the constraint:

minimize 
$$f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$
  
subject to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , (13)

where  $\rho > {\rm 0}$  is the regularization parameter.

#### Definition (Augmented Lagrangian (1969) [6, 7])

The Lagrangian for problem (13) is:

$$\mathcal{L}_{\rho}(\boldsymbol{x},\boldsymbol{\nu}) := f(\boldsymbol{x}) + \boldsymbol{\nu}^{\top} (\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}) + \frac{\rho}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_{2}^{2}.$$
 (14)

This Lagrangian is called the augmented Lagrangian for problem (4).

# Augmented Lagrangian Method (Method of Multipliers)

• The augmented Lagrangian:

$$\mathcal{L}_{
ho}(\boldsymbol{x}, \boldsymbol{
u}) := f(\boldsymbol{x}) + \boldsymbol{
u}^{ op} (\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}) + rac{
ho}{2} \| \boldsymbol{A}\boldsymbol{x} - \boldsymbol{b} \|_2^2.$$

• Recall Eqs. (8) and (9):

$$\begin{split} \mathbf{x}^{(k+1)} &:= \arg\min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}^{(k)}), \\ \mathbf{\nu}^{(k+1)} &:= \boldsymbol{\nu}^{(k)} + \eta^{(k)} (\mathbf{A} \mathbf{x}^{(k+1)} - \mathbf{b}), \end{split}$$

• We can use this augmented Lagrangian in Eqs. (8) and (9):

$$\boldsymbol{x}^{(k+1)} := \arg\min_{\boldsymbol{x}} \mathcal{L}_{\rho}(\boldsymbol{x}, \boldsymbol{\nu}^{(k)}), \tag{15}$$

$$\boldsymbol{\nu}^{(k+1)} := \boldsymbol{\nu}^{(k)} + \rho(\boldsymbol{A}\boldsymbol{x}^{(k+1)} - \boldsymbol{b}), \tag{16}$$

where we use  $\rho$  for the step size of updating the dual variable. This method is called the *augmented Lagrangian method* or the **method of multipliers** (1969) [6, 7, 8].

Alternating Direction Method of Multipliers (ADMM)

# Alternating Direction Method of Multipliers (ADMM)

- Alternating Direction Method of Multipliers (ADMM), proposed in 1976 [9, 10, 11], has been used in many recent machine learning and signal processing papers.
- The usefulness and goal for using ADMM (and other distributed methods) are two-fold:
  - it makes the problem distributed and parallelizable on several servers,
  - it makes it possible to solve an optimization problem with multiple variables.
- Consider the following problem:

$$\begin{array}{l} \underset{x_1, x_2}{\text{minimize}} & f_1(x_1) + f_2(x_2) \\ \text{subject to} & \boldsymbol{A} x_1 + \boldsymbol{B} x_2 = \boldsymbol{c}, \end{array}$$

$$(17)$$

which is an optimization over two variables  $x_1$  and  $x_2$ .

• The augmented Lagrangian for this problem is:

$$\mathcal{L}_{\rho}(\mathbf{x}_{1}, \mathbf{x}_{2}, \boldsymbol{\nu}) = f_{1}(\mathbf{x}_{1}) + f_{2}(\mathbf{x}_{2}) + \boldsymbol{\nu}^{\top}(\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}\|_{2}^{2}.$$
 (18)

 We can alternate between updating the primal variables x<sub>1</sub> and x<sub>2</sub> and the dual variable ν until convergence of these variables:

$$\mathbf{x}_{1}^{(k+1)} := \arg\min_{\mathbf{x}_{1}} \mathcal{L}_{\rho}(\mathbf{x}_{1}, \mathbf{x}_{2}^{(k)}, \boldsymbol{\nu}^{(k)}),$$
(19)

$$\mathbf{x}_{2}^{(k+1)} := \arg\min_{\mathbf{x}_{2}} \mathcal{L}_{\rho}(\mathbf{x}_{1}^{(k+1)}, \mathbf{x}_{2}, \boldsymbol{\nu}^{(k)}),$$
(20)

$$\boldsymbol{\nu}^{(k+1)} := \boldsymbol{\nu}^{(k)} + \rho(\boldsymbol{A}\boldsymbol{x}_1^{(k+1)} + \boldsymbol{B}\boldsymbol{x}_2^{(k+1)} - \boldsymbol{c}). \tag{21}$$

# Alternating Direction Method of Multipliers (ADMM)

- Note that the order of updating primal and dual variables is important and the dual variable should be updated after the primal variables but the order of updating primal variables is not important.
- A good survey/tutorial on ADMM is by Boyd in 2011 [11].
- As was explained before, Eqs. (19) and (20) can be performed by any optimization method such as calculating the gradient of augmented Lagrangian w.r.t. x<sub>1</sub> and x<sub>2</sub>, respectively, and using a few (or even one) iterations of gradient descent for each of these equations.

### Simplifying Equations in ADMM

• The last term in the augmented Lagrangian, Eq. (18), can be restated as:

$$\nu^{\top} (\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}\|_{2}^{2}$$

$$= \nu^{\top} (\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}\|_{2}^{2} + \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2} - \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2}$$

$$= \frac{\rho}{2} \Big( \|\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}\|_{2}^{2} + \frac{1}{\rho^{2}} \|\boldsymbol{\nu}\|_{2}^{2} + \frac{2}{\rho} \nu^{\top} (\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c}) \Big) - \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2}$$

$$\stackrel{(a)}{=} \frac{\rho}{2} \|\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c} + \frac{1}{\rho} \boldsymbol{\nu}\|_{2}^{2} - \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2}$$

$$\stackrel{(b)}{=} \frac{\rho}{2} \|\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c} + \mathbf{u}\|_{2}^{2} - \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2}.$$

where (a) is because of the square of summation of two terms and (b) is because we define  ${m u}:=(1/
ho){m 
u}.$ 

- The last term -(1/(2ρ)) ||ν||<sup>2</sup><sub>2</sub> is constant w.r.t. the primal variables x<sub>1</sub> and x<sub>2</sub> so we can drop that term from Lagrangian when updating the primal variables.
- Hence, the augmented Lagrangian can be restated as:

$$\mathcal{L}_{\rho}(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{u}) = f_{1}(\mathbf{x}_{1}) + f_{2}(\mathbf{x}_{2}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2} - \mathbf{c} + \mathbf{u}\|_{2}^{2} + \text{constant.}$$
(22)

# Simplifying Equations in ADMM

• The augmented Lagrangian:

$$\mathcal{L}_{
ho}(\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{u}) = f_1(\boldsymbol{x}_1) + f_2(\boldsymbol{x}_2) + rac{
ho}{2} \|\boldsymbol{A}\boldsymbol{x}_1 + \boldsymbol{B}\boldsymbol{x}_2 - \boldsymbol{c} + \boldsymbol{u}\|_2^2 + ext{constant.}$$

For updating x<sub>1</sub> and x<sub>2</sub>, the terms f<sub>2</sub>(x<sub>2</sub>) and f(x<sub>1</sub>) are constant, respectively, and can be dropped (because here arg min is important and not the minimum value). Hence, Eqs. (19), (20), and (21) can be restated as:

$$\mathbf{x}_{1}^{(k+1)} := \arg\min_{\mathbf{x}_{1}} \left( f_{1}(\mathbf{x}_{1}) + \frac{\rho}{2} \| \mathbf{A}\mathbf{x}_{1} + \mathbf{B}\mathbf{x}_{2}^{(k)} - \mathbf{c} + \mathbf{u}^{(k)} \|_{2}^{2} \right),$$
(23)

$$\boldsymbol{x}_{2}^{(k+1)} := \arg\min_{\boldsymbol{x}_{2}} \left( f_{2}(\boldsymbol{x}_{2}) + \frac{\rho}{2} \| \boldsymbol{A} \boldsymbol{x}_{1}^{(k+1)} + \boldsymbol{B} \boldsymbol{x}_{2} - \boldsymbol{c} + \boldsymbol{u}^{(k)} \|_{2}^{2} \right),$$
(24)

$$\boldsymbol{u}^{(k+1)} := \boldsymbol{u}^{(k)} + \rho(\boldsymbol{A}\boldsymbol{x}_1^{(k+1)} + \boldsymbol{B}\boldsymbol{x}_2^{(k+1)} - \boldsymbol{c}).$$
<sup>(25)</sup>

- Again, Eqs. (23) and (24) can be performed by one or few steps of gradient descent or any other optimization method.
- The convergence of ADMM for non-convex and non-smooth functions has been analyzed in [12].

- ADMM can be extended to several equality and inequality constraints for several optimization variables [13, 14].
- Consider the following optimization problem with *m* optimization variables and an equality and inequality constraint for every variable:

$$\begin{array}{ll} \underset{\{\mathbf{x}_i\}_{i=1}^m}{\text{minimize}} & \sum_{i=1}^m f_i(\mathbf{x}_i) \\ \text{subject to} & y_i(\mathbf{x}_i) \le 0, \ i \in \{1, \dots, m\}, \\ & h_i(\mathbf{x}_i) = 0, \ i \in \{1, \dots, m\}. \end{array}$$
(26)

• We can convert every inequality constraint to equality constraints by this technique [13, 14]:

$$y_i(\boldsymbol{x}_i) \leq 0 \equiv y'_i(\boldsymbol{x}_i) := (\max(0, y_i(\boldsymbol{x}_i)))^2 = 0.$$

Hence, the problem becomes:

$$\begin{array}{ll} \underset{\{\boldsymbol{x}_i\}_{i=1}^m}{\min i mize} & \sum_{i=1}^m f_i(\boldsymbol{x}_i) \\ \text{subject to} & y_i'(\boldsymbol{x}_i) = 0, \ i \in \{1, \dots, m\}, \\ & h_i(\boldsymbol{x}_i) = 0, \ i \in \{1, \dots, m\}. \end{array}$$

We found:

$$\begin{array}{ll} \underset{\{\mathbf{x}_i\}_{i=1}^m}{\min } & \sum_{i=1}^m f_i(\mathbf{x}_i) \\ \text{subject to} & y_i'(\mathbf{x}_i) = 0, \ i \in \{1, \dots, m\}, \\ & h_i(\mathbf{x}_i) = 0, \ i \in \{1, \dots, m\}. \end{array}$$

• Having dual variables  $\lambda = [\lambda_1, \dots, \lambda_m]^\top$  and  $\nu = [\nu_1, \dots, \nu_m]^\top$  and regularization parameter  $\rho > 0$ , the augmented Lagrangian for this problem is:

$$\mathcal{L}_{\rho}(\{\mathbf{x}_{i}\}_{i=1}^{m}, \boldsymbol{\nu}', \boldsymbol{\nu}) = \sum_{i=1}^{m} f_{i}(\mathbf{x}_{i}) + \sum_{i=1}^{m} \lambda_{i} y_{i}'(\mathbf{x}_{i}) + \sum_{i=1}^{m} \nu_{i} h_{i}(\mathbf{x}_{i}) + \frac{\eta}{2} \sum_{i=1}^{m} (y_{i}'(\mathbf{x}_{i}))^{2} + \frac{\rho}{2} \sum_{i=1}^{m} (h_{i}(\mathbf{x}_{i}))^{2} = \sum_{i=1}^{m} f_{i}(\mathbf{x}_{i}) + \boldsymbol{\lambda}^{\top} \mathbf{y}'(\mathbf{x}) + \boldsymbol{\nu}^{\top} \mathbf{h}(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{y}'(\mathbf{x})\|_{2}^{2} + \frac{\rho}{2} \|\mathbf{h}(\mathbf{x})\|_{2}^{2},$$
(27)

where  $\mathbb{R}^m \ni \mathbf{y}'(\mathbf{x}) := [y'_1(\mathbf{x}_1), \dots, y'_m(\mathbf{x}_m)]^\top$  and  $\mathbb{R}^m \ni \mathbf{h}(\mathbf{x}) := [h_1(\mathbf{x}_1), \dots, h_m(\mathbf{x}_m)]^\top$ .

We found:

$$\mathcal{L}_{\rho}(\{\mathbf{x}_i\}_{i=1}^m, \boldsymbol{\nu}', \boldsymbol{\nu}) = \sum_{i=1}^m f_i(\mathbf{x}_i) + \boldsymbol{\lambda}^\top \mathbf{y}'(\mathbf{x}) + \boldsymbol{\nu}^\top \mathbf{h}(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{y}'(\mathbf{x})\|_2^2 + \frac{\rho}{2} \|\mathbf{h}(\mathbf{x})\|_2^2.$$

Updating the primal and dual variables are performed as [13, 14]:

$$\begin{aligned} \mathbf{x}_{i}^{(k+1)} &:= \arg\min_{\mathbf{x}_{i}} \mathcal{L}_{\rho}(\mathbf{x}_{i}, \lambda_{i}^{(k)}, \nu_{i}^{(k)}), \, \forall i \in \{1, \dots, m\}, \\ \mathbf{\lambda}^{(k+1)} &:= \mathbf{\lambda}^{(k)} + \rho \, \mathbf{y}'(\mathbf{x}^{(k+1)}), \\ \mathbf{\nu}^{(k+1)} &:= \mathbf{\nu}^{(k)} + \rho \, \mathbf{h}(\mathbf{x}^{(k+1)}). \end{aligned}$$

 Note that as the Lagrangian is completely decomposable by the *i* indices, the optimization for every *i*-th primal or dual variable does not depend on other indices; in other words, the terms of other indices become constant for every index.

• The last terms in the augmented Lagrangian, Eq. (27), can be restated as:

$$\begin{split} \lambda^{\top} \mathbf{y}'(\mathbf{x}) + \boldsymbol{\nu}^{\top} \mathbf{h}(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{y}'(\mathbf{x})\|_{2}^{2} + \frac{\rho}{2} \|\mathbf{h}(\mathbf{x})\|_{2}^{2} \\ &= \lambda^{\top} \mathbf{y}'(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{y}'(\mathbf{x})\|_{2}^{2} + \frac{1}{2\rho} \|\lambda\|_{2}^{2} - \frac{1}{2\rho} \|\lambda\|_{2}^{2} \\ &\quad + \boldsymbol{\nu}^{\top} \mathbf{h}(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{h}(\mathbf{x})\|_{2}^{2} + \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2} - \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2} \\ &= \frac{\rho}{2} \Big( \|\mathbf{y}'(\mathbf{x})\|_{2}^{2} + \frac{1}{\rho^{2}} \|\lambda\|_{2}^{2} + \frac{2}{\rho} \lambda^{\top} \mathbf{y}'(\mathbf{x}) \Big) - \frac{1}{2\rho} \|\lambda\|_{2}^{2} \\ &\quad + \frac{\rho}{2} \Big( \|\mathbf{h}(\mathbf{x})\|_{2}^{2} + \frac{1}{\rho^{2}} \|\boldsymbol{\nu}\|_{2}^{2} + \frac{2}{\rho} \boldsymbol{\nu}^{\top} \mathbf{h}(\mathbf{x}) \Big) - \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2} \\ &= \frac{\rho}{2} \|\mathbf{y}'(\mathbf{x}) + \frac{1}{\rho} \lambda\|_{2}^{2} - \frac{1}{2\rho} \|\lambda\|_{2}^{2} + \frac{\rho}{2} \|\mathbf{h}(\mathbf{x}) + \frac{1}{\rho} \boldsymbol{\nu}\|_{2}^{2} - \frac{1}{2\rho} \|\boldsymbol{\nu}\|_{2}^{2} \\ &\stackrel{(a)}{=} \frac{\rho}{2} \|\mathbf{y}'(\mathbf{x}) + \mathbf{u}_{\lambda}\|_{2}^{2} + \frac{\rho}{2} \|\mathbf{h}(\mathbf{x}) + \mathbf{u}_{\nu}\|_{2}^{2} - \text{constant}, \end{split}$$

where (a) is because we define  $oldsymbol{u}_{\lambda}:=(1/
ho)oldsymbol{\lambda}$  and  $oldsymbol{u}_{
u}:=(1/
ho)oldsymbol{
u}.$ 

The augmented Lagrangian was:

$$\mathcal{L}_{\rho}(\{\boldsymbol{x}_i\}_{i=1}^m, \boldsymbol{\nu}', \boldsymbol{\nu}) = \sum_{i=1}^m f_i(\boldsymbol{x}_i) + \boldsymbol{\lambda}^\top \boldsymbol{y}'(\boldsymbol{x}) + \boldsymbol{\nu}^\top \boldsymbol{h}(\boldsymbol{x}) + \frac{\rho}{2} \|\boldsymbol{y}'(\boldsymbol{x})\|_2^2 + \frac{\rho}{2} \|\boldsymbol{h}(\boldsymbol{x})\|_2^2.$$

We simplified the last terms in the augmented Lagrangian:

$$\begin{split} \boldsymbol{\lambda}^{\top} \boldsymbol{y}'(\boldsymbol{x}) + \boldsymbol{\nu}^{\top} \boldsymbol{h}(\boldsymbol{x}) + \frac{\rho}{2} \|\boldsymbol{y}'(\boldsymbol{x})\|_2^2 + \frac{\rho}{2} \|\boldsymbol{h}(\boldsymbol{x})\|_2^2 \\ &= \frac{\rho}{2} \|\boldsymbol{y}'(\boldsymbol{x}) + \boldsymbol{u}_{\lambda}\|_2^2 + \frac{\rho}{2} \|\boldsymbol{h}(\boldsymbol{x}) + \boldsymbol{u}_{\nu}\|_2^2 - \text{constant.} \end{split}$$

• Hence, the Lagrangian can be restated as:

$$\begin{split} \mathcal{L}_{\rho}(\{\bm{x}_{i}\}_{i=1}^{m}, \bm{u}_{\lambda}, \bm{u}_{\nu}) &= \sum_{i=1}^{m} f_{i}(\bm{x}_{i}) + \frac{\rho}{2} \|\bm{y}'(\bm{x}) + \bm{u}_{\lambda}\|_{2}^{2} + \frac{\rho}{2} \|\bm{h}(\bm{x}) + \bm{u}_{\nu}\|_{2}^{2} + \text{constant} \\ &= \sum_{i=1}^{m} f_{i}(\bm{x}_{i}) + \frac{\rho}{2} \sum_{i=1}^{m} \left[ (y_{i}'(\bm{x}_{i}) + u_{\lambda,i})^{2} + (h_{i}(\bm{x}_{i}) + u_{\nu,i})^{2} \right] + \text{constant}, \end{split}$$

where  $u_{\lambda,i} = (1/\rho)\lambda_i$  and  $u_{\nu,i} = (1/\rho)\nu_i$  are the *i*-th elements of  $\boldsymbol{u}_{\lambda}$  and  $\boldsymbol{u}_{\nu}$ , respectively.

• The Lagrangian was restated as:

$$\mathcal{L}_{\rho}(\{\mathbf{x}_{i}\}_{i=1}^{m}, \mathbf{u}_{\lambda}, \mathbf{u}_{\nu}) = \sum_{i=1}^{m} f_{i}(\mathbf{x}_{i}) + \frac{\rho}{2} \|\mathbf{y}'(\mathbf{x}) + \mathbf{u}_{\lambda}\|_{2}^{2} + \frac{\rho}{2} \|\mathbf{h}(\mathbf{x}) + \mathbf{u}_{\nu}\|_{2}^{2} + \text{constant}$$
$$= \sum_{i=1}^{m} f_{i}(\mathbf{x}_{i}) + \frac{\rho}{2} \sum_{i=1}^{m} [(y_{i}'(\mathbf{x}_{i}) + u_{\lambda,i})^{2} + (h_{i}(\mathbf{x}_{i}) + u_{\nu,i})^{2}] + \text{constant}.$$

Hence, updating variables can be restated as:

$$\mathbf{x}_{i}^{(k+1)} := \arg\min_{\mathbf{x}_{i}} \left( f_{i}(\mathbf{x}_{i}) + \frac{\rho}{2} \left[ (y_{i}'(\mathbf{x}_{i}) + u_{\lambda,i}^{(k)})^{2} + (h_{i}(\mathbf{x}_{i}) + u_{\nu,i}^{(k)})^{2} \right] \right), \forall i \in \{1, \dots, m\},$$
(28)

$$u_{\lambda,i}^{(k+1)} := u_{\lambda,i}^{(k)} + \rho \, y_i'(\mathbf{x}_i^{(k+1)}), \, \forall i \in \{1, \dots, m\}$$
<sup>(29)</sup>

$$u_{\nu,i}^{(k+1)} := u_{\nu,i}^{(k)} + \rho h_i(\boldsymbol{x}_i^{(k+1)}), \, \forall i \in \{1, \dots, m\}.$$
(30)

Use of ADMM for Distributed Optimization

#### Use of ADMM for Distributed Optimization

- ADMM is one of the most well-known algorithms for distributed optimization.
- If the problem can be divided into several disjoint blocks (i.e., several primal variables), we can solve the optimization for each primal variable on a separate core or server (see Eq. (28) for every *i*):

$$\mathbf{x}_{i}^{(k+1)} := \arg\min_{\mathbf{x}_{i}} \left( f_{i}(\mathbf{x}_{i}) + \frac{\rho}{2} \left[ (y_{i}'(\mathbf{x}_{i}) + u_{\lambda,i}^{(k)})^{2} + (h_{i}(\mathbf{x}_{i}) + u_{\nu,i}^{(k)})^{2} \right] \right), \forall i \in \{1, \dots, m\}.$$

Hence, in every iteration of ADMM, the update of primal variables can be performed in parallel by distributed servers.

• At the end of each iteration, the updated primal variables are gathered in a central server so that the update of dual variable(s) is performed (see Eqs. (29) and (30)):

$$u_{\lambda,i}^{(k+1)} := u_{\lambda,i}^{(k)} + \rho y_i'(\mathbf{x}_i^{(k+1)}), \, \forall i \in \{1, \dots, m\}, \\ u_{\nu,i}^{(k+1)} := u_{\nu,i}^{(k)} + \rho h_i(\mathbf{x}_i^{(k+1)}), \, \forall i \in \{1, \dots, m\}.$$

- Then, the updated dual variable(s) is sent to the distributed servers so they update their primal variables. This procedure is repeated until convergence of primal and dual variables.
- In this sense, ADMM is performed similar to the approach of federated learning [15, 16].

- We can convert a non-distributed optimization problem to a distributed optimization problem to solve it using ADMM. Many recent machine learning and signal processing papers are using this technique.
- Univariate optimization problem: Consider a regular non-distributed problem with one optimization variable *x*:

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & \sum_{i=1}^{m} f_i(\mathbf{x}) \\ \text{subject to} & y_i(\mathbf{x}) \leq 0, \ i \in \{1, \dots, m\}, \\ & h_i(\mathbf{x}) = 0, \ i \in \{1, \dots, m\}. \end{array}$$

$$(31)$$

This problem can be stated as:

$$\begin{array}{ll} \underset{\{x_i\}_{i=1}^m}{\mininize} & \sum_{i=1}^m f_i(x_i) \\ \text{subject to} & y_i(x_i) \le 0, \ i \in \{1, \dots, m\}, \\ & h_i(x_i) = 0, \ i \in \{1, \dots, m\}, \\ & x_i = z, \ i \in \{1, \dots, m\}, \end{array}$$
(32)

where we introduce *m* variables  $\{\mathbf{x}_i\}_{i=1}^m$  and use the trick  $\mathbf{x}_i = \mathbf{z}, \forall i$  to make them equal to one variable.

- Eq. (32) is similar to Eq. (26) except that it has 2*m* equality constraints rather than *m* equality constraints.
- Hence, we can use ADMM updates similarly to Eqs. (28), (29), and (30) but with slight change because of the additional *m* constraints.
- We introduce *m* new dual variables for constraints x<sub>i</sub> = z, ∀i and update those dual variables as well as other variables. The augmented Lagrangian also has some additional terms for the new constraints.
- The Lagrangian and ADMM updates of this are not stated here because of its similarity to the previous equations.
- This is a good technique to make a problem distributed, use ADMM for solving it, and solving it in parallel servers.

• Multivariate optimization problem: Consider a regular non-distributed problem with multiple optimization variables  $\{x_i\}_{i=1}^m$ :

$$\begin{array}{l} \underset{\{\mathbf{x}\}_{i=1}^{m}}{\text{minimize}} & \sum_{i=1}^{m} f_i(\mathbf{x}_i) \\ \text{subject to} & \mathbf{x}_i \in \mathcal{S}_i, \ i \in \{1, \dots, m\}, \end{array}$$

$$(33)$$

where  $\mathbf{x}_i \in S_i$  can be any constraint such as belonging to a set  $S_i$ , an equality constraint, or an inequality constraint.

• We can embed the constraint in the objective function using an indicator function:

$$\underset{\{\boldsymbol{x}\}_{i=1}^{m}}{\text{minimize}} \quad \sum_{i=1}^{m} (f_{i}(\boldsymbol{x}_{i}) + \phi_{i}(\boldsymbol{x}_{i})),$$

where  $\phi_i(\mathbf{x}_i) := \mathbb{I}(\mathbf{x}_i \in S_i)$  is zero if  $\mathbf{x}_i \in S_i$  and is infinity otherwise.

s

This problem can be stated as:

$$\min_{\substack{\{\boldsymbol{x}_i\}_{i=1}^m \\ \text{subject to}}} \sum_{i=1}^m \left( f_i(\boldsymbol{x}_i) + \phi_i(\boldsymbol{z}_i) \right)$$

$$\text{subject to} \quad \boldsymbol{x}_i = \boldsymbol{z}_i, \ i \in \{1, \dots, m\},$$

$$(34)$$

where we introduce a variable  $z_i$  for every  $x_i$ , use the introduced variable for the second term in the objective function, and we equate them in the constraint.

We found:

$$\min_{\{\boldsymbol{x}_i\}_{i=1}^m} \sum_{i=1}^m \left(f_i(\boldsymbol{x}_i) + \phi_i(\boldsymbol{z}_i)\right)$$

subject to  $\mathbf{x}_i = \mathbf{z}_i, i \in \{1, \ldots, m\}.$ 

As the constraints x<sub>i</sub> − z<sub>i</sub> = 0, ∀i are equality constraints, we can use Eqs. (23), (24), and (25) as ADMM updates for this problem:

$$\mathbf{x}_{i}^{(k+1)} := \arg\min_{\mathbf{x}_{i}} \left( f_{i}(\mathbf{x}_{i}) + \frac{\rho}{2} \| \mathbf{x}_{i} - \mathbf{z}_{i}^{(k)} + \mathbf{u}_{i}^{(k)} \|_{2}^{2} \right), \quad \forall i \in \{1, \dots, m\},$$
(35)

$$\mathbf{z}_{i}^{(k+1)} := \arg\min_{\mathbf{z}_{i}} \left( \phi_{i}(\mathbf{z}_{i}) + \frac{\rho}{2} \| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i} + \mathbf{u}_{i}^{(k)} \|_{2}^{2} \right), \quad \forall i \in \{1, \dots, m\},$$
(36)

$$\boldsymbol{u}_{i}^{(k+1)} := \boldsymbol{u}_{i}^{(k)} + \rho(\boldsymbol{x}_{i}^{(k+1)} + \boldsymbol{z}_{i}^{(k+1)}), \, \forall i \in \{1, \ldots, m\}.$$

• Comparing Eqs. (35) and (36) with the proximal operator,  $\operatorname{prox}_{\lambda g}(\mathbf{x}) := \arg \min_{\mathbf{u}} \left( g(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{x}\|_2^2 \right)$ , shows that these ADMM updates can be written as proximal mappings:

$$\begin{aligned} \mathbf{x}_{i}^{(k+1)} &:= \operatorname{prox}_{\frac{1}{\rho}\rho_{i}}(\mathbf{z}_{i}^{(k)} - \mathbf{u}_{i}^{(k)}), \,\forall i \in \{1, \dots, m\}, \\ \mathbf{z}_{i}^{(k+1)} &:= \operatorname{prox}_{\frac{1}{\rho}\rho_{i}}(\mathbf{x}_{i}^{(k+1)} + \mathbf{u}_{i}^{(k)}), \,\forall i \in \{1, \dots, m\}, \\ \mathbf{u}_{i}^{(k+1)} &:= \mathbf{u}_{i}^{(k)} + \rho(\mathbf{x}_{i}^{(k+1)} + \mathbf{z}_{i}^{(k+1)}), \,\forall i \in \{1, \dots, m\}, \end{aligned}$$
that  $\|\mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}\| + u_{i}^{(k)}\|^{2} = \|\mathbf{z}_{i} - \mathbf{x}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k)}\|^{2}$ 

if we notice that  $\|\boldsymbol{x}_i^{(k+1)} - \boldsymbol{z}_i + \boldsymbol{u}_i^{(k)}\|_2^2 = \|\boldsymbol{z}_i - \boldsymbol{x}_i^{(k+1)} - \boldsymbol{u}_i^{(k)}\|_2^2.$ 

- Note that in many papers, such as [17], we only have m = 1. In that case, we only have two primal variables x and z.
- According to the lemma we had (that the proximal operator of indicator function is projection), as the function  $\phi_i(.)$  is an indicator function, Eq. (37):

$$m{z}_i^{(k+1)} := {\sf prox}_{rac{1}{
ho} \phi_i}(m{x}_i^{(k+1)} + m{u}_i^{(k)}), \, orall i \in \{1, \dots, m\},$$

can be implemented by projection onto the set  $S_i$ :

$$\boldsymbol{z}_{i}^{(k+1)} := \Pi_{\mathcal{S}_{i}}(\boldsymbol{x}_{i}^{(k+1)} + \boldsymbol{u}_{i}^{(k)}), \, \forall i \in \{1, \ldots, m\}.$$

As an example, assume the variables are all matrices so we have X<sub>i</sub>, Z<sub>i</sub>, and U<sub>i</sub>. If the set S<sub>i</sub> is the cone of orthogonal matrices, the constraint X<sub>i</sub> ∈ S<sub>i</sub> would be X<sub>i</sub><sup>T</sup>X<sub>i</sub> = I. In this case, the update of matrix variable Z<sub>i</sub> would be done by setting the singular values of (x<sub>i</sub><sup>(k+1)</sup> + u<sub>i</sub><sup>(k)</sup>) to one (recall projection onto the cone of the orthogonal matrices).

- This example is in my own paper (Image Structural Component Analysis) [18].
- Let image be divided into b blocks. Our goal is to find a p-dimensional subspace for reconstruction of every image block (each block as q pixels and we have p ≪ q).



• Considering all the *b* blocks in an image, the problem is:

$$\begin{array}{ll} \underset{\boldsymbol{U}_i \in \mathbb{R}^{q \times p}}{\text{minimize}} & \sum_{i=1}^{b} || \check{\mathbf{x}}_i - \boldsymbol{U}_i \boldsymbol{U}_i^\top \check{\mathbf{x}}_i ||_S, \\ \text{subject to} & \boldsymbol{U}_i^\top \boldsymbol{U}_i = \boldsymbol{I}, \quad \forall i \in \{1, \dots, b\}, \end{array}$$
(38)

where  $\check{\mathbf{x}}_i \in \mathbb{R}^q$  and  $U_i \in \mathbb{R}^{q \times p}$  are the *i*-th block and the bases of its subspace, respectively.

• We had:

$$\begin{array}{ll} \underset{U_i \in \mathbb{R}^{q \times p}}{\text{minimize}} & \sum_{i=1}^{b} || \breve{\mathbf{x}}_i - \boldsymbol{U}_i \boldsymbol{U}_i^\top \breve{\mathbf{x}}_i ||_{\mathcal{S}}, \\ \text{subject to} & \boldsymbol{U}_i^\top \boldsymbol{U}_i = \boldsymbol{I}, \quad \forall i \in \{1, \dots, b\}, \end{array}$$

• We convert it to:

$$\begin{array}{ll} \underset{\boldsymbol{U}_{i},\boldsymbol{V}_{i}\in\mathbb{R}^{q\times p}}{\text{minimize}} & \sum_{i=1}^{b} \left(f(\boldsymbol{U}_{i})+h(\boldsymbol{V}_{i})\right), \\ \text{subject to} & \boldsymbol{U}-\boldsymbol{V}=\boldsymbol{0}, \end{array}$$
(39)

---**j**--- - - -,

where  $f(\boldsymbol{U}_i) := ||\boldsymbol{x}_i - \boldsymbol{U}_i \boldsymbol{U}_i^\top \boldsymbol{x}_i||_S$  and  $h(\boldsymbol{V}_i) := \mathbb{I}(\boldsymbol{V}_i^\top \boldsymbol{V}_i = \boldsymbol{I}).$ 

• The (squared) SSIM distance, which we denote by  $||.||_S$ , is [19]:

$$\mathbb{R} \ni ||\check{\mathbf{x}}_1 - \check{\mathbf{x}}_2||_S := 1 - \mathsf{SSIM}(\check{\mathbf{x}}_1, \check{\mathbf{x}}_2) = \frac{||\check{\mathbf{x}}_1 - \check{\mathbf{x}}_2||_2^2}{||\check{\mathbf{x}}_1||_2^2 + ||\check{\mathbf{x}}_2||_2^2 + c},\tag{40}$$

• The I(.) denotes the indicator function which is zero if its condition is satisfied and is infinite otherwise.

• The **U** and **V** are defined as union of partitions to form an image-form array, i.e.,  $\mathbf{U} := \bigcup_{i=1}^{b} \mathbf{U}_i$  and  $\mathbf{V} := \bigcup_{i=1}^{b} \mathbf{V}_i$  [17].

• Eq. (39) was:

$$\min_{\boldsymbol{U}_i, \boldsymbol{V}_i \in \mathbb{R}^{q \times p}} \sum_{i=1}^{b} (f(\boldsymbol{U}_i) + h(\boldsymbol{V}_i)),$$

subject to  $\boldsymbol{U} - \boldsymbol{V} = \boldsymbol{0}.$ 

• The augmented Lagrangian for Eq. (39) is:

$$\begin{aligned} \mathcal{L}_{\rho} &= \sum_{i=1}^{b} \left( f(\boldsymbol{U}_{i}) + h(\boldsymbol{V}_{i}) \right) + \operatorname{tr} (\boldsymbol{\Lambda}^{\top} (\boldsymbol{U} - \boldsymbol{V})) + (\rho/2) ||\boldsymbol{U} - \boldsymbol{V}||_{F}^{2} \\ &= \sum_{i=1}^{b} \left( f(\boldsymbol{U}_{i}) + h(\boldsymbol{V}_{i}) \right) + (\rho/2) ||\boldsymbol{U} - \boldsymbol{V} + \boldsymbol{J}||_{F}^{2} - (\rho/2) ||\boldsymbol{\Lambda}||_{F}^{2}, \end{aligned}$$

where  $\|.\|_F$  is the Frobenius norm,  $\mathbf{\Lambda} := \bigcup_{i=1}^{b} \mathbf{\Lambda}_i$  is the Lagrange multiplier,  $\rho > 0$  is a parameter, and  $\mathbf{J} := (1/\rho)\mathbf{\Lambda} = (1/\rho) \bigcup_{i=1}^{b} \mathbf{\Lambda}_i = \bigcup_{i=1}^{b} \mathbf{J}_i$ .

• Note that the term  $(\rho/2) ||\mathbf{A}||_{\mathbf{F}}^2$  is a constant with respect to  $\mathbf{U}$  and  $\mathbf{V}$  and can be dropped. The updates of  $\mathbf{U}, \mathbf{V}$ , and  $\mathbf{J}$  are done as [11, 17]:

$$\boldsymbol{U}_{i}^{(k+1)} := \arg\min_{\boldsymbol{U}_{i}} \left( f(\boldsymbol{U}_{i}) + (\rho/2) || \boldsymbol{U}_{i} - \boldsymbol{V}_{i}^{(k)} + \boldsymbol{J}_{i}^{(k)} ||_{F}^{2} \right),$$
(41)

$$\boldsymbol{V}_{i}^{(k+1)} := \arg\min_{\boldsymbol{V}_{i}} \left( h(\boldsymbol{V}_{i}) + (\rho/2) || \boldsymbol{U}_{i}^{(k+1)} - \boldsymbol{V}_{i} + \boldsymbol{J}_{i}^{(k)} ||_{F}^{2} \right),$$
(42)

$$J^{(k+1)} := J^{(k)} + U^{(k+1)} - V^{(k+1)}.$$
(43)



**Fig. 1.** Examples from the training dataset: (a) original image, (b) contrast stretched, (c) Gaussian noise, (d) luminance enhanced, (e) Gaussian blurring, (f) salt & pepper impulse noise, and (g) JPEG distortion.



Fig. 2. The first dimension of the trained (a) U, (b) V, and (c) J for ISCA.

# Acknowledgement

- Some slides of this slide deck are inspired by the lectures of Prof. Stephen Boyd at the Stanford University.
- Our tutorial also has the materials of this slide deck: [20]

#### References

- Q. Li, Z. Zhu, and G. Tang, "Alternating minimizations converge to second-order optimal solutions," in *International Conference on Machine Learning*, pp. 3935–3943, 2019.
- [2] P. Jain and P. Kar, "Non-convex optimization for machine learning," *arXiv preprint arXiv:1712.07897*, 2017.
- [3] G. B. Dantzig and P. Wolfe, "Decomposition principle for linear programs," Operations research, vol. 8, no. 1, pp. 101–111, 1960.
- [4] J. F. Benders, "Partitioning procedures for solving mixed-variables programming problems," *Numerische mathematik*, vol. 4, no. 1, pp. 238–252, 1962.
- [5] H. Everett III, "Generalized Lagrange multiplier method for solving problems of optimum allocation of resources," *Operations research*, vol. 11, no. 3, pp. 399–417, 1963.
- [6] M. R. Hestenes, "Multiplier and gradient methods," Journal of optimization theory and applications, vol. 4, no. 5, pp. 303–320, 1969.
- [7] M. J. Powell, "A method for nonlinear constraints in minimization problems," *Optimization*, pp. 283–298, 1969.
- [8] D. P. Bertsekas, "The method of multipliers for equality constrained problems," *Constrained optimization and Lagrange multiplier methods*, pp. 96–157, 1982.

# References (cont.)

- [9] D. Gabay and B. Mercier, "A dual algorithm for the solution of nonlinear variational problems via finite element approximation," *Computers & mathematics with applications*, vol. 2, no. 1, pp. 17–40, 1976.
- [10] R. Glowinski and A. Marrocco, "Finite element approximation and iterative methods of solution for 2-D nonlinear magnetostatic problems," in *Proceeding of International Conference on the Computation of Electromagnetic Fields (COMPUMAG)*, 1976.
- S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.
- [12] Y. Wang, W. Yin, and J. Zeng, "Global convergence of ADMM in nonconvex nonsmooth optimization," *Journal of Scientific Computing*, vol. 78, no. 1, pp. 29–63, 2019.
- [13] J. Giesen and S. Laue, "Distributed convex optimization with many convex constraints," arXiv preprint arXiv:1610.02967, 2016.
- [14] J. Giesen and S. Laue, "Combining ADMM and the augmented Lagrangian method for efficiently handling many constraints," in *International Joint Conference on Artificial Intelligence*, pp. 4525–4531, 2019.
- [15] J. Konečný, B. McMahan, and D. Ramage, "Federated optimization: Distributed optimization beyond the datacenter," arXiv preprint arXiv:1511.03575, 2015.

# References (cont.)

- [16] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [17] D. Otero, D. La Torre, O. V. Michailovich, and E. R. Vrscay, "Alternate direction method of multipliers for unconstrained structural similarity-based optimization," in *International Conference Image Analysis and Recognition*, pp. 20–29, Springer, 2018.
- [18] B. Ghojogh, F. Karray, and M. Crowley, "Principal component analysis using structural similarity index for images," in *Image Analysis and Recognition: 16th International Conference, ICIAR 2019, Waterloo, ON, Canada, August 27–29, 2019, Proceedings, Part I* 16, pp. 77–88, Springer, 2019.
- [19] D. Brunet, E. R. Vrscay, and Z. Wang, "On the mathematical properties of the structural similarity index," *IEEE Transactions on Image Processing*, vol. 21, no. 4, pp. 1488–1499, 2012.
- [20] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "KKT conditions, first-order and second-order optimization, and distributed optimization: Tutorial and survey," arXiv preprint arXiv:2110.01858, 2021.