

Tutorial on Distributed Optimization

Amal Feriani

University of Manitoba

December 17, 2023

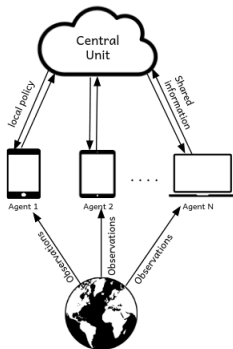
Table of Contents

- 1 Introduction and Motivation
- 2 Decomposition
- 3 Dual Decomposition
- 4 Augmented Lagrangian Methods/ Methods of Multipliers
- 5 Alternating Direction Method of Multipliers
- 6 Consensus
- 7 Conclusions

Table of Contents

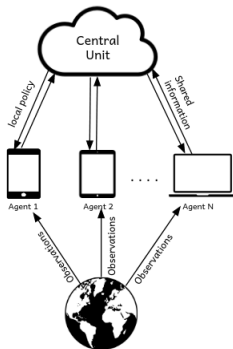
- 1 Introduction and Motivation
- 2 Decomposition
- 3 Dual Decomposition
- 4 Augmented Lagrangian Methods/ Methods of Multipliers
- 5 Alternating Direction Method of Multipliers
- 6 Consensus
- 7 Conclusions

Motivation

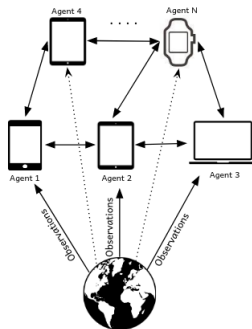


Centralized

Motivation

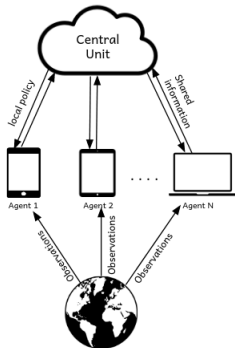


Centralized

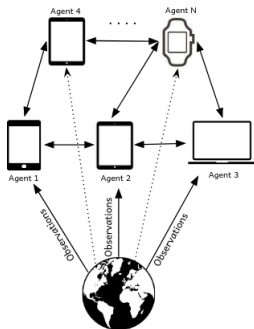


Networked Agents

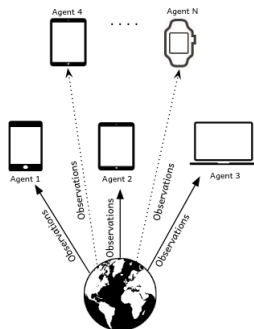
Motivation



Centralized



Networked Agents



Fully Decentralized

- Fully decentralized learning scheme

- Fully decentralized learning scheme
 - Non-stationarity problem
 - Local solutions may not necessarily lead to global ones
 - Non-cooperative setting

- Fully decentralized learning scheme
 - Non-stationarity problem
 - Local solutions may not necessarily lead to global ones
 - Non-cooperative setting
- Centralized learning scheme

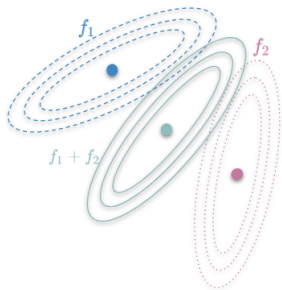
- Fully decentralized learning scheme
 - Non-stationarity problem
 - Local solutions may not necessarily lead to global ones
 - Non-cooperative setting
- Centralized learning scheme
 - Existence of a central coordinator

- Fully decentralized learning scheme
 - Non-stationarity problem
 - Local solutions may not necessarily lead to global ones
 - Non-cooperative setting
- Centralized learning scheme
 - Existence of a central coordinator
- Networked Optimization scheme

- Fully decentralized learning scheme
 - Non-stationarity problem
 - Local solutions may not necessarily lead to global ones
 - Non-cooperative setting
- Centralized learning scheme
 - Existence of a central coordinator
- Networked Optimization scheme
 - Existence of a communication structure to enable coordination
 - Information exchange with neighboring nodes

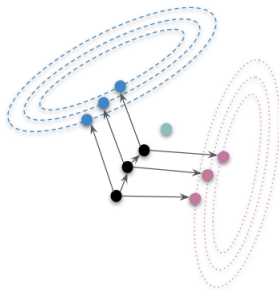
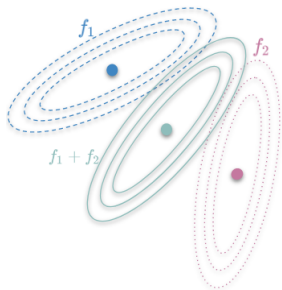
Motivation

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f^i(\mathbf{x}) \quad (1)$$



Motivation

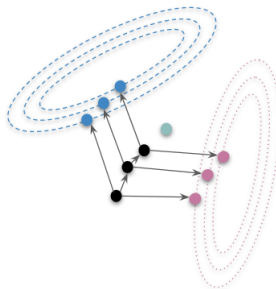
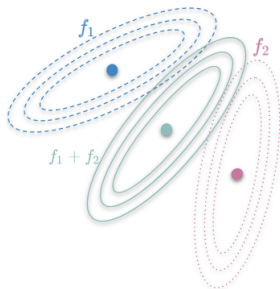
$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f^i(\mathbf{x}) \quad (1)$$



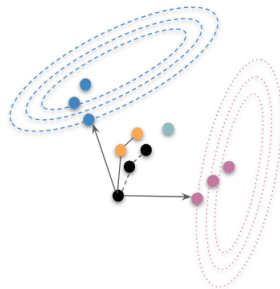
$$x_k = x_{k-1} - \alpha_k \sum_{i=0}^N g^i(x_{k-1})$$

Motivation

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f^i(\mathbf{x}) \quad (1)$$



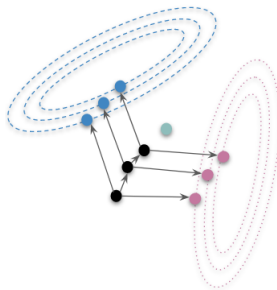
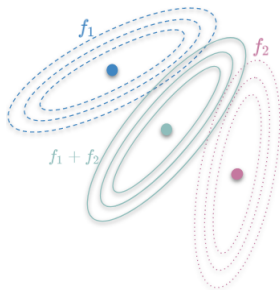
$$x_k = x_{k-1} - \alpha_k \sum_{i=0}^N g^i(x_{k-1})$$



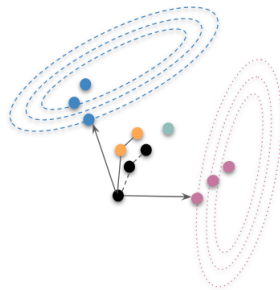
$$x_k^i = x_{k-1}^i - \alpha_k^i g^i(x_{k-1}^i)$$

Motivation

$$\mathbf{f}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f^i(\mathbf{x}) \quad (2)$$



slow (only 1 update)



Potentially faster but may not converge

Table of Contents

- 1 Introduction and Motivation
- 2 Decomposition**
- 3 Dual Decomposition
- 4 Augmented Lagrangian Methods/ Methods of Multipliers
- 5 Alternating Direction Method of Multipliers
- 6 Consensus
- 7 Conclusions

Decomposable Problems

- Optimization problems with *coupling* variables

$$\max_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f^i(\mathbf{x})$$

Decomposable Problems

- Optimization problems with *coupling* variables

$$\max_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f^i(\mathbf{x})$$

- Optimization problems with *coupling* constraints

$$\max_{x_i} \frac{1}{N} \sum_{i=1}^N f^i(x_i) \text{ s.t } C(\mathbf{x}) = B$$

Power allocation in wireless communication

$$\begin{aligned} \max_{p_i} \quad & \sum_{i=0}^N \log\left(1 + \frac{g_i p_i}{\sigma^2}\right) \\ \text{s.t.} \quad & \sum_{i=0}^N p_i = P_{\max} \leftarrow \text{coupling constraint} \end{aligned}$$

Decomposition: Motivation

Decomposition benefits for MAS:

- Computational Scalability and flexibility
- Tractability
- Privacy
-

Table of Contents

- 1 Introduction and Motivation
- 2 Decomposition
- 3 Dual Decomposition**
- 4 Augmented Lagrangian Methods/ Methods of Multipliers
- 5 Alternating Direction Method of Multipliers
- 6 Consensus
- 7 Conclusions

Dual Ascent Algorithm

Consider the problem

$$\min_{\mathbf{x}} f(\mathbf{x}) \text{ s.t. } A\mathbf{x} = \mathbf{b}$$

- **Lagrangian**

$$L(\mathbf{x}, \mathbf{u}) = f(\mathbf{x}) + (A\mathbf{x} - \mathbf{b})^T \mathbf{u}$$

- **Dual function**

$$g(\mathbf{u}) = \min_{\mathbf{x}} L(\mathbf{x}, \mathbf{u})$$

- **Conjugate function/Fenchel conjugate**

$$\begin{aligned} f^*(\mathbf{y}) &= \max_{\mathbf{x}} [\mathbf{y}^T \mathbf{x} - f(\mathbf{x})] \\ -f^*(\mathbf{y}) &= \min_{\mathbf{x}} [f(\mathbf{x}) - \mathbf{y}^T \mathbf{x}] \\ g(\mathbf{u}) &= -f^*(-A^T \mathbf{u}) - \mathbf{b}^T \mathbf{u} \end{aligned}$$

Reminder: Conjugate functions

- If f is closed and convex, $(f^*)^* = f$ and

$$x \in \partial f^* \iff y \in \partial f \iff x \in \arg \min_z f(z) - y^T z \quad (3)$$

- If f is strictly convex, $\nabla f^*(y) = \arg \min_z f(z) - y^T z$
- f^* is always convex even if f is not convex (i.e affine in y)

Dual Ascent Algorithm

Dual Problem:

$$\max_{\mathbf{u}} g(\mathbf{u}) = \min_{\mathbf{u}} f^*(-A^T \mathbf{u}) + \mathbf{b}^T \mathbf{u}$$

$$\implies \partial g(\mathbf{u}) = A \partial f^*(-A^T \mathbf{u}) - \mathbf{b}$$

$$\implies \partial g(\mathbf{u}) = A \mathbf{x}^* - \mathbf{b} \text{ where } \mathbf{x}^* \in \arg \min_{\mathbf{x}} f(\mathbf{x}) + \mathbf{u}^T A \mathbf{x} \quad (3)$$

Apply sub/gradient methods to the dual : $u_{k+1} = u_k + t_k \partial g(u_k)$

Dual Ascent Algorithm

Dual Problem:

$$\max_{\mathbf{u}} g(\mathbf{u}) = \min_{\mathbf{x}} f^*(-A^T \mathbf{u}) + \mathbf{b}^T \mathbf{u}$$

$$\implies \partial g(\mathbf{u}) = A \partial f^*(-A^T \mathbf{u}) - \mathbf{b}$$

$$\implies \partial g(\mathbf{u}) = A\mathbf{x}^* - \mathbf{b} \text{ where } \mathbf{x}^* \in \arg \min_{\mathbf{x}} f(\mathbf{x}) + \mathbf{u}^T A\mathbf{x} \quad (3)$$

Apply sub/gradient methods to the dual : $u_{k+1} = u_k + t_k \partial g(u_k)$

Algorithm 2: Dual Sub-gradient Method

Start with an initial dual λ_0 and repeat for K iterations; t_0 initial learning rate

for $k = 1, \dots, K$ **do**

$\mathbf{x}_k \in \arg \min_{\mathbf{x}} f(\mathbf{x}) + \mathbf{u}_{k-1}^T A\mathbf{x}$
 $\mathbf{u}_k = \mathbf{u}_{k-1} + t_k(A\mathbf{x}_k - \mathbf{b})$

end

Separable Problem Example

Primal Problem

$$\begin{aligned} \min_{x_1, x_2} \quad & f_1(x_1) + f_2(x_2) \\ \text{s.t.} \quad & x_1 + x_2 \leq P \end{aligned}$$

Dual Problem

$$\begin{aligned} \max_u \quad & -f_1^*(-u) - f_2^*(-u) - Pu \\ \text{s.t.} \quad & u \geq 0 \end{aligned}$$

The Lagrangian is separable in f :

$$L(x_1, x_2, u) = f_1(x_1) + ux_1 + f_2(x_2) + ux_2 - Pu$$

$$g(u) = g_1(u) + g_2(u) - Pu; \quad g_i(u) = \min_{x_i} f_i(x_i) + ux_i = -f_i^*(-u), i = 1, 2$$

Applying Dual Ascent will result in

$$u_{k+1} = u_k + t_k \text{Proj}_{u \geq 0} \left\{ x_k^1 + x_k^2 - b \right\} = u_k + t_k \max \left\{ x_k^1 + x_k^2 - b, 0 \right\}$$

$$x_k^i = \arg \min_{x^i} f_i(x^i) + u_{k-1} x^i, \quad i = 1, 2 \implies \text{decoupled}$$

Dual Decomposition

More generally,

$$\begin{array}{ll} \min_{\mathbf{x}_i} \sum_{i=1}^B f_i(\mathbf{x}_i) & \text{separable} \\ \text{s.t } A\mathbf{x} = \mathbf{b} & \text{not separable} \end{array}$$

Dual Decomposition

More generally,

$$\begin{array}{ll} \min_{\mathbf{x}_i} \sum_{i=1}^B f_i(\mathbf{x}_i) & \text{separable} \\ \text{s.t } A\mathbf{x} = \mathbf{b} & \text{not separable} \end{array}$$

Algorithm 4: Dual Decomposition Method (coupled constraints)

Start with an initial dual λ_0 and repeat for K iterations; t_0 initial learning rate **for** $k = 1, \dots, K$ **do**

$$\begin{array}{l} \mathbf{x}_k^i = \arg \min_{\mathbf{x}^i} f_i(\mathbf{x}^i) + \mathbf{u}_{k-1}^T A^i \mathbf{x}^i \\ \mathbf{u}_k = \mathbf{u}_{k-1} + t_k (\sum_{i=1}^B A^i \mathbf{x}_k^i - \mathbf{b}) \end{array}$$

end

- The minimization problems w.r.t to x_i are independent
- Coordination is done through adjusting u to maximize $g(u)$

Dual Decomposition for coupling variables

Primal Problem

$$\min_{x_1, x_2, y} f_1(x_1, y) + f_2(x_2, y)$$

Transform the problem to

$$\begin{aligned} \min_{x_1, x_2, y_1, y_2} f_1(x_1, y_1) + f_2(x_2, y_2) & : y_1, y_2 \text{ are } \mathbf{local\ copies} \text{ of } y \\ \text{s.t } y_1 = y_2 & : \mathbf{consensus\ constraint} \end{aligned}$$

Algorithm 5: Dual Decomposition for coupling variables

Start with an initial dual λ_0 and repeat for K iterations; t_0 initial learning rate

for $k = 1 \dots K$ **do**

$$\begin{aligned} & x_k^i, y_k^i = \arg \min_{x^i, y^i} f_i(x^i, y^i) + u_{k-1} y^i, \quad i = 1, 2 \\ & u_k = u_{k-1} + t_k(y_k^1 - y_k^2) \end{aligned}$$

end

But, the convergence to a feasible iterates (i.e $y_2 - y_1 \neq 0$) is not always guaranteed.

Discussion

- The dual variable u is equivalent to a price of a shared resource
- If the resource is over-utilized (i.e $P - x_1 - x_2 < 0$), the price is increased and decreased otherwise.
- But, the price never gets negative !
- A centralized unit is needed for coordination to update u

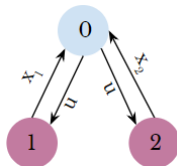


Figure: Price Coordination interpretation of DD (see slide 16)

- **Subgradient** methods are slow, step size selection is difficult
- **Gradient** methods are applicable if the dual function is differentiable (which is not always the case, requires strict convexity of the primal function)
- The convergence guarantees of the Dual Ascent method do not necessarily guarantee the convergence of the primal (i.e Slater's Theorem)

Problem 1

$$\begin{aligned} \min_{x_1, x_2} & 2|x_1 - 2| + 4|x_2 - 4| \\ \text{s.t. } & x_1 + x_2 = 5 \\ & 0 \leq x_1 \leq 10, 0 \leq x_2 \leq 10 \end{aligned}$$

The optimal solution : $x_1^* = 1, x_2^* = 4$ with optimal value is $f^* = 2$

Problem 2

$$\begin{aligned} \min_{x_1, x_2} & 2(x_1 - 2)^2 + 4(x_2 - 4)^2 \\ \text{s.t. } & x_1 + x_2 = 5 \\ & 0 \leq x_1 \leq 10, 0 \leq x_2 \leq 10 \end{aligned}$$

The solution for this problem is $x_1^* = 4/3, x_2^* = 11/3$ and the optimal value is $f^* = 4/3$.

Discussion

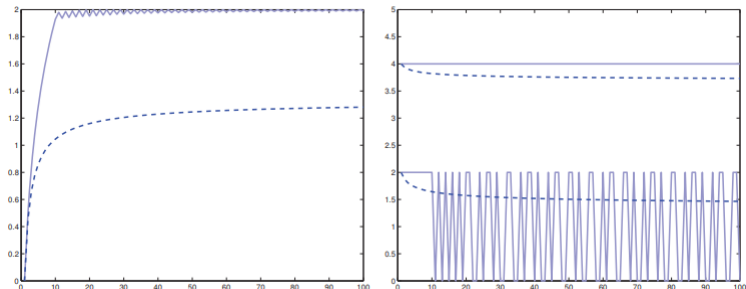


Figure: When the dual function is non-smooth, the dual objective (left, full lines) converges while the primal iterates (right, full lines) do not. When the dual function is differentiable, both objective and iterates converge asymptotically (dashed lines) [1]

Table of Contents

- 1 Introduction and Motivation
- 2 Decomposition
- 3 Dual Decomposition
- 4 Augmented Lagrangian Methods/ Methods of Multipliers**
- 5 Alternating Direction Method of Multipliers
- 6 Consensus
- 7 Conclusions

Augmented Lagrangian Methods (ALM)/ Methods of Multipliers

Adds a penalty term to the primal objective to ensure the smoothness of the dual function

Augmented Lagrangian Methods (ALM)/ Methods of Multipliers

Adds a penalty term to the primal objective to ensure the smoothness of the dual function

Considers the modified problem, for a parameter $\rho > 0$

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \\ \text{s.t } \mathbf{Ax} = \mathbf{b} \end{aligned} \tag{4}$$

Augmented Lagrangian Methods (ALM)/ Methods of Multipliers

Adds a penalty term to the primal objective to ensure the smoothness of the dual function

Considers the modified problem, for a parameter $\rho > 0$

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \\ \text{s.t } \mathbf{Ax} = \mathbf{b} \end{aligned} \quad (4)$$

Uses the **augmented Lagrangian**

$$L_{\rho}(\mathbf{x}, \mathbf{u}) = f(\mathbf{x}) + \mathbf{u}^T (\mathbf{Ax} - \mathbf{b}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \quad (5)$$

Augmented Lagrangian Methods (ALM)/ Methods of Multipliers

Adds a penalty term to the primal objective to ensure the smoothness of the dual function

Considers the modified problem, for a parameter $\rho > 0$

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) + \frac{\rho}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ \text{s.t } A\mathbf{x} = \mathbf{b} \end{aligned} \quad (4)$$

Uses the **augmented Lagrangian**

$$L_{\rho}(\mathbf{x}, \mathbf{u}) = f(\mathbf{x}) + \mathbf{u}^T (A\mathbf{x} - \mathbf{b}) + \frac{\rho}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 \quad (5)$$

Repeats, $k = 1 \dots K$

$$\begin{aligned} \mathbf{x}_k &= \arg \min_{\mathbf{x}} f(\mathbf{x}) + \mathbf{u}_{k-1}^T A\mathbf{x} + \frac{\rho}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ \mathbf{u}_k &= \mathbf{u}_{k-1} + \rho(A\mathbf{x}_k - \mathbf{b}) \end{aligned}$$

ALM convergence

- Optimality conditions

Primal feasibility $Ax^* - b = 0$; stationarity $\partial f(x^*) + A^T u^* = 0$

- x_k minimizes $L_\rho(x, y_k)$

$$\begin{aligned} \mathbf{x}_k &= \arg \min_{\mathbf{x}} f(\mathbf{x}) + \lambda_{k-1}^T A\mathbf{x} + \frac{\rho}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ \implies 0 &\in \partial f(\mathbf{x}_k) + A^T \underbrace{(\lambda_{k-1} + \rho(A\mathbf{x}_k - \mathbf{b}))}_{\lambda_k} \end{aligned}$$

- the dual update makes the pair (x_k, y_k) satisfy the **stationarity** condition
- Under mild conditions, $A\mathbf{x}_k - \mathbf{b} \rightarrow 0$ as $k \rightarrow \infty$, KKT conditions are satisfied in the limit
- x_k and u_k converge to solutions

- The objective function stays the same
- The primal objective because strongly convex if A is full rank (f is a convex function)
- The Dual function becomes smooth \implies gradient ascent
- Better convergence guarantees compared to Dual Ascent (f can be non-differentiable, take on value $+\infty$)
- The primal objective is no longer decomposable because of the quadratic term

Proximal Operator

$$\text{prox}_{\rho f}(\tilde{x}) = \arg \min_x f(x) + \frac{1}{2\rho} \|x - \tilde{x}\|^2$$

Proximal Operator

$$\text{prox}_{\rho f}(\tilde{x}) = \arg \min_x f(x) + \frac{1}{2\rho} \|x - \tilde{x}\|^2$$

Proximal Point Algorithm (PPA)[2]

$$x_{t+1} = \text{prox}_{\rho f}(x_t) = (I + \rho \partial f)^{-1}(x_t)$$

Augmented Lagrangian and Proximal Point Methods

$$\textbf{Primal} : \min_x f(x) + g(Ax)$$

$$\textbf{Dual} : \max_u -f^*(-A^T u) - g^*(u) = \max_u -H(u)$$

Applying the PPA in the dual problem:

$$\begin{aligned} u_{t+1} &= \text{prox}_{\rho H}(u_t) \\ &= \arg \min_u f^*(-A^T u) + g^*(u) + \frac{1}{2\rho} \|u - u_t\|_2^2 \end{aligned}$$

Is equivalent to

$$\begin{aligned} u_{t+1} &= u_t + \rho(A\tilde{x} - \tilde{y}) \text{ where} \\ (\tilde{x}, \tilde{y}) &= \arg \min_{x,y} (f(x) + g(y) + u_t^T(Ax - y) + \frac{\rho}{2} \|Ax - y\|_2^2) \end{aligned}$$

Hence, the augmented Lagrangian in the primal is the same as PPA in the dual

Example

$$\min_{\mathbf{x}} \sum_{i=1}^N f_i(\mathbf{x})$$

$$\begin{cases} \min_{\mathbf{x}_i} \sum_{i=1}^N f_i(\mathbf{x}_i) \\ \text{s.t } \mathbf{x}_i = \mathbf{x} \quad i = 1 \dots N \end{cases}$$

$g = \mathbb{1}_{\mathbf{x}}$ (indicator function), $A = I_N$

The augmented Lagrangian method

$$(x_{k+1}^i, x_{k+1}) = \arg \min_{x_i, x} (f_i(x_i) + u_t^T (x_i - x) + \frac{\rho}{2} \|x_i - x\|_2^2)$$

$$u_{k+1}^i = u_k^i + \rho(x_{k+1}^i - x_{k+1})$$

Coupling between \mathbf{x} and \mathbf{x}_i

Table of Contents

- 1 Introduction and Motivation
- 2 Decomposition
- 3 Dual Decomposition
- 4 Augmented Lagrangian Methods/ Methods of Multipliers
- 5 Alternating Direction Method of Multipliers**
- 6 Consensus
- 7 Conclusions

Alternating Direction Method of Multipliers (ADMM)

Combines the best of both methods (decomposability and convergence guarantees) [3]

Alternating Direction Method of Multipliers (ADMM)

Combines the best of both methods (decomposability and convergence guarantees) [3]

Consider a problem of the form:

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{y}} \quad & f(\mathbf{x}) + g(\mathbf{y}) \\ \text{s.t.} \quad & A\mathbf{x} + B\mathbf{y} = \mathbf{c} \end{aligned} \tag{6}$$

Two blocks of variables with separable objectives

Alternating Direction Method of Multipliers (ADMM)

Combines the best of both methods (decomposability and convergence guarantees) [3]

Consider a problem of the form:

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{y}} \quad & f(\mathbf{x}) + g(\mathbf{y}) \\ \text{s.t.} \quad & A\mathbf{x} + B\mathbf{y} = \mathbf{c} \end{aligned} \tag{6}$$

Two blocks of variables with separable objectives

The augmented Lagrangian is

$$L_{\rho}(x, y, u) = f(x) + g(y) + u^T(Ax + By - c) + \frac{\rho}{2} \|Ax + By - c\|_2^2$$

Alternating Direction Method of Multipliers (ADMM)

Combines the best of both methods (decomposability and convergence guarantees) [3]

Consider a problem of the form:

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{y}} & f(\mathbf{x}) + g(\mathbf{y}) \\ \text{s.t. } & \mathbf{Ax} + \mathbf{By} = \mathbf{c} \end{aligned} \tag{6}$$

Two blocks of variables with separable objectives

The augmented Lagrangian is

$$L_{\rho}(x, y, u) = f(x) + g(y) + u^T (Ax + By - c) + \frac{\rho}{2} \|Ax + By - c\|_2^2$$

Applying the method of multipliers, we need to find (x_k, y_k) by solving the joint minimization:

$$(x_k, y_k) = \arg \min_{x, y} L_{\rho}(x, y, u_{k-1})$$

Alternating Direction Method of Multipliers (ADMM)

ADMM splits the joint minimization into two steps :

$$x_k = \arg \min_x L_\rho(x, y_{k-1}, u_{k-1}) \quad (7)$$

$$y_k = \arg \min_y L_\rho(x_k, y, u_{k-1}) \quad (8)$$

$$u_k = u_{k-1} + \rho(Ax_k + By_k - c) \quad (9)$$

⇒ Splitting since we minimize over x with y fixed, and vice versa
(**Gauss-Seidel method**)

Example

$$\begin{cases} \min_{\mathbf{x}_i} \sum_{i=1}^N f_i(\mathbf{x}_i) \\ \text{s.t } \mathbf{x}_i = \mathbf{x} \quad i = 1 \dots N \end{cases}$$

The augmented Lagrangian

$$L_{\rho}(x, y, u) = \sum_{i=1}^N f_i(x_i) + u^T (x_i - x) + \frac{\rho}{2} \|x_i - x\|_2^2$$

ADMM iterates

$$x_{k+1}^i = \arg \min_{x_i} (f_i(x_i) + (u_k^i)^T x_i + \frac{\rho}{2} \|x_i - x_k\|_2^2)$$

$$x_{k+1} = \frac{1}{N} \sum_{i=1}^N x_{k+1}^i + \frac{1}{\rho} u_k^i$$

$$u_{k+1}^i = u_k^i + \rho (x_{k+1}^i - x_{k+1})$$

Convergence guarantees

Under modest assumptions on f , g , the ADMM iterates satisfy, for any $\rho > 0$ [4]:

- **Residual convergence/Primal feasibility:** $r_k = Ax_k + By_k - c \rightarrow 0$ as $k \rightarrow \infty$
- **Primal convergence:** $f(x_k) + g(y_k) \rightarrow f^* + g^*$
- **Dual convergence:** $u_k \rightarrow u^*$

Assumption 1. The (extended-real-valued) functions $f : \mathbb{R}^n \mapsto \mathbb{R} \cup \{+\infty\}$ and $g : \mathbb{R}^m \mapsto \mathbb{R} \cup \{+\infty\}$ are closed, proper, and convex. $\implies \text{epi} f = \{(x, t) \in \mathbb{R}^n \times \mathbb{R} \mid f(x) \leq t\}$ is a closed nonempty convex set.

Assumption 2. The unaugmented Lagrangian L has a saddle point.

Note: No further assumptions on A and B

Convergence rate: Theory still being developed but ADMM has similar convergence rate as 1st order methods (linear) (i.e [5] linear convergence rate when one of the functions is strongly convex)

Scaled form ADMM

Let the residual $r = Ax + Bz - c$ and the scaled dual variable $w = \frac{u}{\rho}$

$$u^T r + \frac{\rho}{2} \|r\|_2^2 = \frac{\rho}{2} \|r + \frac{1}{\rho} u\|_2^2 - \frac{2}{2\rho} \|u\|_2^2 = \frac{\rho}{2} \|r + w\|_2^2 - \frac{\rho}{2} \|w\|_2^2$$

The augmented Lagrangian becomes

$$L_\rho(x, y, w) = f(x) + g(y) + \frac{\rho}{2} \|r + w\|_2^2 - \frac{\rho}{2} \|w\|_2^2$$

And the ADMM iterates

$$x_{k+1} = \arg \min_x f(x) + \frac{\rho}{2} \|Ax + By_k - c + w_k\|_2^2$$

$$y_{k+1} = \arg \min_y g(y) + \frac{\rho}{2} \|Ax_{k+1} + By - c + w_k\|_2^2$$

$$w_{k+1} = w_k + Ax_{k+1} + By_{k+1} - c$$

The residual at iteration k , $r_k = Ax_k + By_k - c$,

$$w_{k+1} = w_0 + \sum_{i=1}^k r_i \longrightarrow \text{running sum of residuals}$$

Douglas–Rachford algorithm

Consider the problem, where f and g are closed convex functions

$$\min_x f(x) + g(x)$$

Douglas–Rachford iteration [6]:

$$y_{k+1} = y_k + \text{prox}_g(2x_k - y_k) - x_k$$

$$x_{k+1} = \text{prox}_f(y_{k+1})$$

An equivalent formulation

$$u_{k+1} = \text{prox}_g(2x_k - y_k)$$

$$x_{k+1} = \text{prox}_f(y_k + u_{k+1} - x_k)$$

$$y_{k+1} = y_k + u_{k+1} - x_{k+1}$$

Let $w_k = x_k - y_k$

$$u_{k+1} = \text{prox}_g(x_k + w_k)$$

$$x_{k+1} = \text{prox}_f(u_{k+1} - w_k)$$

$$w_{k+1} = w_k + x_{k+1} - u_{k+1}$$

Douglas–Rachford algorithm and ADMM

ADMM is equivalent to applying Douglas–Rachford in the Dual :

$$\min_{x,y} f(x) + g(y) \text{ s.t } Ax + By = c$$

Dual Problem: $\max_z -c^T z - f^*(-A^T z) - g^*(-B^T z)$

Apply the Douglas–Rachford method to minimize

$$\underbrace{c^T z + f^*(-A^T z)}_{f_1(z)} + \underbrace{g^*(-B^T z)}_{f_2(z)}$$

$$u_{k+1} = \text{prox}_{f_1}(z_k + w_k) \quad (10)$$

$$z_{k+1} = \text{prox}_{f_2}(u_{k+1} - w_k) \quad (11)$$

$$w_{k+1} = w_k + z_{k+1} - u_{k+1} \quad (12)$$

Recall that the PPA algorithm in the dual is equivalent to the augmented Lagrangian in the primal

Douglas–Rachford algorithm and ADMM

(10) \implies

$$\tilde{x} = \arg \min_x f(x) + (z_k + w_k)^T (Ax - b) + \rho/2 \|Ax - b\|_2^2 \quad (13)$$

$$u_{k+1} = (z_k + w_k) + \rho(A\tilde{x} - b)$$

(11) \implies

$$\tilde{y} = \arg \min_y f(y) + (z_k)^T (By) + \rho/2 \|A\tilde{x} + By - b\|_2^2 \quad (14)$$

$$z_{k+1} = z_k + \rho(A\tilde{x} + B\tilde{y} - b) \quad (15)$$

From (13), (14) and (15), we can derive the ADMM iterates

- Modest accuracy in a handful of iterations, but it requires a large number of iterations for a highly accurate solution (like a first-order method)
- Choice of ρ can greatly influence practical convergence of ADMM:
 - ρ too large \rightarrow not enough emphasis on minimizing $f + g$
 - ρ too small \rightarrow not enough emphasis on feasibility
 - Vector Approximate Message Passing (VAMP) find the optimal ρ value under the Gaussian approximation of the messages
- Transforming a problem to ADMM format can be subtle and leads to different algorithms

ADMM for constrained convex optimization

Consider problem

$$\min_x f(x) \text{ s.t } x \in C$$

Transform the problem to ADMM format

$$\min_{x,y} f(x) + g(y) \text{ s.t } x - y = 0$$

where $g = \mathbb{I}_C$

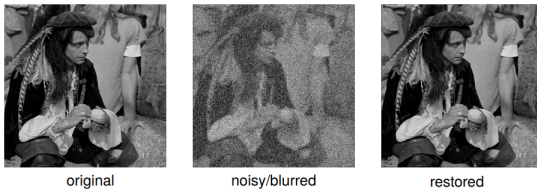
$$x_{k+1} = \arg \min_x (f(x) + \frac{\rho}{2} \|x - y_k + u_k\|_2^2)$$

$$y_{k+1} = \arg \min_y (g(y) + \frac{\rho}{2} \|x_{k+1} + u_k - y\|_2^2) = \Pi_C(x_{k+1} + u_k)$$

$$u_{k+1} = u_k + x_{k+1} - y_{k+1}$$

ADMM for computer vision

Image Deblurring[7]



Sparse + Low rank matrix decomposition[8]

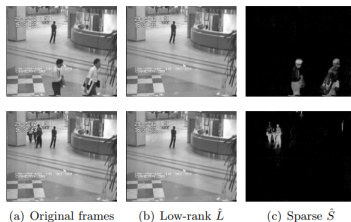


Table of Contents

- 1 Introduction and Motivation
- 2 Decomposition
- 3 Dual Decomposition
- 4 Augmented Lagrangian Methods/ Methods of Multipliers
- 5 Alternating Direction Method of Multipliers
- 6 Consensus**
- 7 Conclusions

Consensus ADMM

Consider the problem

$$\min_x \sum_{i=1}^N f_i(x)$$

The **Consensus ADMM** starts by transforming the problem to

$$\min_{x_i, x} \sum_{i=1}^N f_i(x_i) \text{ s.t. } x_i = x, \quad i = 1 \dots N$$

Thus, the **decomposable** ADMM iterates

$$x_k^i = \arg \min_{x_i} f_i(x_i) + \rho/2 \|x_i - x_{k-1} + w_{k-1}^i\|_2^2 \quad i = 1, \dots, N$$

$$x_k = 1/N \sum_{i=1}^N (x_k^i + w_{k-1}^i)$$

$$w_k^i = w_{k-1}^i + x_k^i - x_k, \quad i = 1, \dots, N$$

Consensus ADMM

By taking $x_k = \bar{x} = 1/N \sum_{i=0}^N x_k^i$, the previous iterates can be simplified to

$$x_k^i = \arg \min_{x_i} f_i(x_i) + \rho/2 \|x_i - \bar{x}_{k-1} + w_{k-1}^i\|_2^2 \quad i = 1, \dots, N$$

$$w_k^i = w_{k-1}^i + x_k^i - \bar{x}_k, \quad i = 1, \dots, N$$

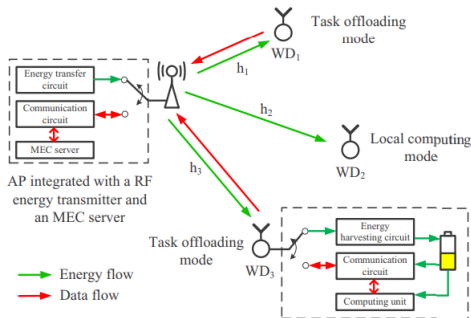
- the $x_i, i = 1, \dots, B$ updates are done **in parallel**
- Consensus ADMM is **communication inefficient** since we need to gather all the x_k^i to update x_k and thus w_k^i
Recently, Group ADMM (GADMM) [9] is proposed to tackle this problem, where a worker i only communicates to two neighbors

$$\min_{x_i} \sum_{i=0}^N f_i(x_i) \text{ s.t } x_i = x_{i+1}, \quad i = 1 \dots N$$

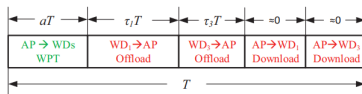
Task Offloading in MEC [10]

- MEC network with an AP and N wireless devices (WD)
- AP broadcasts RF energy to distributed WDs and the WDs harvest that energy for aT time, $a \in [0, 1]$ in each time frame
- Within each time frame, we assume that each WD needs to accomplish a certain computing task (pollution level) based on its local data (measurements)
- WD decides to compute the task locally (mode 0) or offload it to the AP (mode 1) (binary offloading)
- $M = M_0 \cup M_1 = \{1 \dots N\}$ the set of WDs
- Harvested energy $E_i = \mu P h_i a T, i = 1 \dots N$, P is RF energy transmit power of the AP, $\mu \in [0, 1]$ is the energy harvesting efficiency, h_i is the channel gain between the AP and the i -th WD
- Transmission time $\tau_i T, \tau \in [0, 1]$

Task Offloading in MEC



An example 3-user wireless powered MEC system with binary computation offloading



An example time allocation in the 3-user wireless powered MEC network. Only WD1 and WD3 selecting mode 1 offload the task to and download the computation results from the AP

- **Local computation model**

$$r_i = \nu_1 \left(\frac{h_i}{k_i} \right)^{1/3} a^{1/3} \quad (16)$$

where r_i is the local computation rate (bits/s), k_i denotes the computation energy efficiency coefficient of the processor's chip and ν_1 is a fixed parameter

- **offloading mode**

$$r_i = \frac{B\tau_i}{\nu_u} \log_2 \left(1 + \frac{\mu P a h_i^2}{\tau_i N_0} \right) \forall i \in M_1 \quad (17)$$

where ν_u indicates the communication overhead in task offloading, B the bandwidth and N_0 the noise power

Problem formulation

$$(P1) : \underset{\mathcal{M}_0, a, \tau}{\text{maximize}} \quad \sum_{i \in \mathcal{M}_0} w_i \eta_1 \left(\frac{h_i}{k_i} \right)^{\frac{1}{3}} a^{\frac{1}{3}} \quad (7a)$$

$$+ \sum_{j \in \mathcal{M}_1} w_j \varepsilon \tau_j \ln \left(1 + \frac{\eta_2 h_j^2 a}{\tau_j} \right) \quad (7b)$$

$$\text{subject to} \quad \sum_{j \in \mathcal{M}_1} \tau_j + a \leq 1, \quad (7c)$$

$$a \geq 0, \tau_j \geq 0, \forall j \in \mathcal{M}_1, \quad (7d)$$

$$\mathcal{M}_0 \subseteq \mathcal{M}, \mathcal{M}_1 = \mathcal{M} \setminus \mathcal{M}_0. \quad (7e)$$

ADMM formulation

$$\begin{aligned} & \underset{a, \mathbf{z}, \mathbf{x}, \boldsymbol{\tau}, \mathbf{m}}{\text{maximize}} && \sum_{i=1}^N w_i \left\{ (1 - m_i) \eta_1 \left(\frac{h_i}{k_i} \right)^{\frac{1}{3}} x_i^{\frac{1}{3}} \right. \\ & && \left. + m_i \varepsilon \tau_i \ln \left(1 + \frac{\eta_2 h_i^2 x_i}{\tau_i} \right) \right\} \\ & \text{subject to} && \sum_{i=1}^N z_i + a \leq 1, \\ & && x_i = a, z_i = \tau_i \quad i = 1, \dots, N, \\ & && a, z_i, x_i, \tau_i \geq 0, \quad m_i \in \{0, 1\}, \quad i = 1, \dots, N. \end{aligned} \tag{8}$$

Reformulate (P1) as an equivalent integer programming problem by introducing binary decision variables m_i and additional artificial variables x_i and z_i

ADMM formulation

$$\underset{a, \mathbf{z}, \mathbf{x}, \boldsymbol{\tau}, \mathbf{m}}{\text{maximize}} \quad \sum_{i=1}^N q_i(x_i, \tau_i, m_i) + g(\mathbf{z}, a) \quad (9a)$$

$$\text{subject to} \quad x_i = a, \tau_i = z_i \quad i = 1, \dots, N, \quad (9b)$$

$$x_i, \tau_i \geq 0, \quad m_i \in \{0, 1\}, \quad i = 1, \dots, N, \quad (9c)$$

where

$$q_i(x_i, \tau_i, m_i) = w_i \left\{ (1 - m_i) \eta_1 \left(\frac{h_i}{k_i} \right)^{\frac{1}{3}} x_i^{\frac{1}{3}} + m_i \varepsilon \tau_i \ln \left(1 + \frac{\eta_2 h_i^2 x_i}{\tau_i} \right) \right\},$$

and

$$g(\mathbf{z}, a) = \begin{cases} 0, & \text{if } (\mathbf{z}, a) \in \mathcal{G}, \\ -\infty, & \text{otherwise,} \end{cases} \quad (10)$$

where

$$\mathcal{G} = \left\{ (\mathbf{z}, a) \mid \sum_{i=1}^N z_i + a \leq 1, a \geq 0, z_i \geq 0, i = 1, \dots, N \right\}.$$

Notice that $i \in M_0$ can be removed from the constraint set

ADMM formulation

$$\begin{aligned} L(\mathbf{u}, \mathbf{v}, \boldsymbol{\theta}) &= \sum_{i=1}^N q_i(\mathbf{u}) + g(\mathbf{v}) + \sum_{i=1}^N \beta_i (x_i - a) \\ &+ \sum_{i=1}^N \gamma_i (\tau_i - z_i) - \frac{c}{2} \sum_{i=1}^N (x_i - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i - z_i)^2 \end{aligned}$$

ADMM formulation

$$\begin{aligned} L(\mathbf{u}, \mathbf{v}, \boldsymbol{\theta}) = & \sum_{i=1}^N q_i(\mathbf{u}) + g(\mathbf{v}) + \sum_{i=1}^N \beta_i (x_i - a) \\ & + \sum_{i=1}^N \gamma_i (\tau_i - z_i) - \frac{c}{2} \sum_{i=1}^N (x_i - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i - z_i)^2 \end{aligned}$$

Step 1: maximize w.r.t to $u = \{x, \tau, m\}$

$$u_{t+1} \arg \max_u L(u, v_t, \theta_t) \quad (18)$$

$$L(\mathbf{u}, \mathbf{v}, \boldsymbol{\theta}) = \sum_{i=1}^N q_i(\mathbf{u}) + g(\mathbf{v}) + \sum_{i=1}^N \beta_i (x_i - a) \\ + \sum_{i=1}^N \gamma_i (\tau_i - z_i) - \frac{c}{2} \sum_{i=1}^N (x_i - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i - z_i)^2$$

Step 1: maximize w.r.t to $u = \{x, \tau, m\}$

$$u_{t+1} \arg \max_u L(u, v_t, \theta_t) \quad (18)$$

$$\begin{cases} \underset{x_i, \tau_i \geq 0}{\text{maximize}} & w_i \eta_1 \left(\frac{h_i}{k_i} \right)^{\frac{1}{3}} x_i^{\frac{1}{3}} + \beta_i^l x_i + \gamma_i^l \tau_i - \frac{c}{2} (x_i - a^l)^2 \\ & - \frac{c}{2} (\tau_i - z_i^l)^2, \text{ if } m_i = 0, \\ \underset{x_i, \tau_i \geq 0}{\text{maximize}} & w_i \varepsilon \tau_i \ln \left(1 + \frac{\eta_2 h_i^2 x_i}{\tau_i} \right) + \beta_i^l x_i + \gamma_i^l \tau_i \\ & - \frac{c}{2} (x_i - a^l)^2 - \frac{c}{2} (\tau_i - z_i^l)^2, \text{ if } m_i = 1. \end{cases}$$

ADMM Algorithm

$$\begin{aligned} L(\mathbf{u}, \mathbf{v}, \boldsymbol{\theta}) &= \sum_{i=1}^N q_i(\mathbf{u}) + g(\mathbf{v}) + \sum_{i=1}^N \beta_i (x_i - a) \\ &+ \sum_{i=1}^N \gamma_i (\tau_i - z_i) - \frac{c}{2} \sum_{i=1}^N (x_i - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i - z_i)^2 \end{aligned}$$

ADMM Algorithm

$$\begin{aligned} L(\mathbf{u}, \mathbf{v}, \boldsymbol{\theta}) = & \sum_{i=1}^N q_i(\mathbf{u}) + g(\mathbf{v}) + \sum_{i=1}^N \beta_i (x_i - a) \\ & + \sum_{i=1}^N \gamma_i (\tau_i - z_i) - \frac{c}{2} \sum_{i=1}^N (x_i - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i - z_i)^2 \end{aligned}$$

Step 1: maximize w.r.t to $u = \{x, \tau, m\}$

$$u_{t+1} \arg \max_u L(u, v_t, \theta_t) \quad (19)$$

ADMM Algorithm

$$\begin{aligned} L(\mathbf{u}, \mathbf{v}, \boldsymbol{\theta}) = & \sum_{i=1}^N q_i(\mathbf{u}) + g(\mathbf{v}) + \sum_{i=1}^N \beta_i (x_i - a) \\ & + \sum_{i=1}^N \gamma_i (\tau_i - z_i) - \frac{c}{2} \sum_{i=1}^N (x_i - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i - z_i)^2 \end{aligned}$$

Step 1: maximize w.r.t to $u = \{x, \tau, m\}$

$$u_{t+1} \arg \max_u L(u, v_t, \theta_t) \quad (19)$$

$$\begin{cases} \underset{x_i, \tau_i \geq 0}{\text{maximize}} & w_i \eta_1 \left(\frac{h_i}{k_i} \right)^{\frac{1}{3}} x_i^{\frac{1}{3}} + \beta_i^l x_i + \gamma_i^l \tau_i - \frac{c}{2} (x_i - a^l)^2 \\ & - \frac{c}{2} (\tau_i - z_i^l)^2, \text{ if } m_i = 0, \\ \underset{x_i, \tau_i \geq 0}{\text{maximize}} & w_i \varepsilon \tau_i \ln \left(1 + \frac{\eta_2 h_i^2 x_i}{\tau_i} \right) + \beta_i^l x_i + \gamma_i^l \tau_i \\ & - \frac{c}{2} (x_i - a^l)^2 - \frac{c}{2} (\tau_i - z_i^l)^2, \text{ if } m_i = 1. \end{cases}$$

ADMM Algorithm

Step 2: maximize w.r.t to $v = \{z, a\}$

$$u_{t+1} \arg \max_u L(u_{t+1}, v, \theta_t) \quad (20)$$

ADMM Algorithm

Step 2: maximize w.r.t to $\mathbf{v} = \{z, a\}$

$$u_{t+1} \arg \max_u L(u_{t+1}, \mathbf{v}, \theta_t) \quad (20)$$

$$\mathbf{v}^{l+1} =$$

$$\arg \max_{\mathbf{z}, a} \sum_{i=1}^N \beta_i^l (x_i^{l+1} - a) + \sum_{i=1}^N \gamma_i^l (\tau_i^{l+1} - z_i) \\ - \frac{c}{2} \sum_{i=1}^N (x_i^{l+1} - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i^{l+1} - z_i)^2$$

$$\text{subject to} \quad \sum_{i=1}^N z_i + a \leq 1, \quad a \geq 0, \quad z_i \geq 0, i = 1, \dots, N.$$

ADMM Algorithm

Step 2: maximize w.r.t to $\mathbf{v} = \{z, a\}$

$$u_{t+1} \arg \max_u L(u_{t+1}, \mathbf{v}, \theta_t) \quad (20)$$

$$\mathbf{v}^{l+1} =$$

$$\arg \max_{\mathbf{z}, a} \sum_{i=1}^N \beta_i^l (x_i^{l+1} - a) + \sum_{i=1}^N \gamma_i^l (\tau_i^{l+1} - z_i) \\ - \frac{c}{2} \sum_{i=1}^N (x_i^{l+1} - a)^2 - \frac{c}{2} \sum_{i=1}^N (\tau_i^{l+1} - z_i)^2$$

$$\text{subject to} \quad \sum_{i=1}^N z_i + a \leq 1, \quad a \geq 0, \quad z_i \geq 0, \quad i = 1, \dots, N.$$

Step 3: minimize w.r.t to $\theta = \{\beta, \gamma\}$

$$\beta_i^{l+1} = \beta_i^l - c(x_i^{l+1} - a^{l+1}), \quad i = 1, \dots, N,$$

$$\gamma_i^{l+1} = \gamma_i^l - c(\tau_i^{l+1} - z_i^{l+1}), \quad i = 1, \dots, N.$$

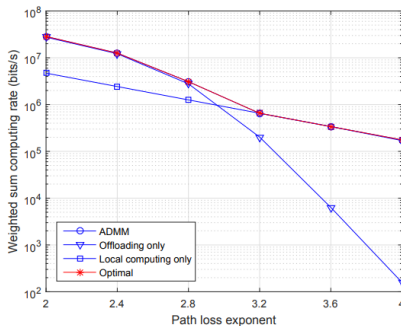
Baselines

- **Optimal**: exhaustively enumerates all the 2^N combinations of N WDs' computing modes;
- **Offloading only**: all the WDs offload their tasks to the AP, $M_0 = \{\}$;
- **Local computing only**: all the WDs perform computations locally, $M_0 = M$.

Results I

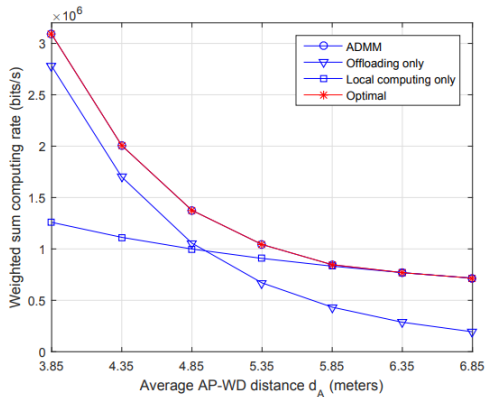
Baselines

- **Optimal**: exhaustively enumerates all the 2^N combinations of N WDs' computing modes;
- **Offloading only**: all the WDs offload their tasks to the AP, $M_0 = \{\}$;
- **Local computing only**: all the WDs perform computations locally, $M_0 = M$.



Under different path loss exponents ($N=10$)

Results II



Under different average AP-WD distance (N=10)

Results III

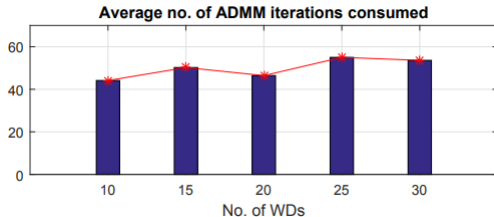
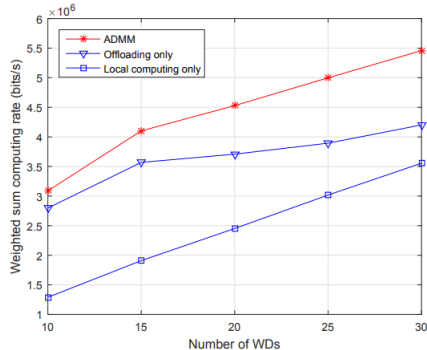


Table of Contents

- 1 Introduction and Motivation
- 2 Decomposition
- 3 Dual Decomposition
- 4 Augmented Lagrangian Methods/ Methods of Multipliers
- 5 Alternating Direction Method of Multipliers
- 6 Consensus
- 7 Conclusions**

- Decomposition is necessary for decentralized optimization
- Coordination between sub-problems can be costly and communication inefficient
- ADMM offers good convergence guarantees with mild assumptions on the objective functions
- ADMM may result in different algorithms when the problem is formulated differently

What we didn't cover

- More on consensus algorithms
- Recent work on distributed gradient algorithms

- ECE236C - Optimization Methods for Large-Scale Systems, Prof. L. Vandenberghe, UCLA
- EE364b - Convex Optimization II, Prof. S. Boyd, Stanford
- 10-725: Convex Optimization, Prof. R. Tibshirani, CMU

- [1] B. Yang and M. Johansson, *Distributed Optimization and Games: A Tutorial Overview*, vol. 406, pp. 109–148. 10 2010.
- [2] R. T. Rockafellar, “Monotone operators and the proximal point algorithm,” *SIAM journal on control and optimization*, vol. 14, no. 5, pp. 877–898, 1976.
- [3] R. Glowinski and A. Marroco, “Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de dirichlet non linéaires,” *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique*, vol. 9, no. R2, pp. 41–76, 1975.

References II

- [4] S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. 2011.
- [5] R. Nishihara, L. Lessard, B. Recht, A. Packard, and M. Jordan, “A general analysis of the convergence of admm,” in *International Conference on Machine Learning*, pp. 343–352, PMLR, 2015.
- [6] J. Eckstein and D. P. Bertsekas, “On the douglas—rachford splitting method and the proximal point algorithm for maximal monotone operators,” *Mathematical Programming*, vol. 55, no. 1-3, pp. 293–318, 1992.
- [7] D. O’Connor and L. Vandenberghe, “Primal-dual decomposition by operator splitting and applications to image deblurring,” *SIAM Journal on Imaging Sciences*, vol. 7, no. 3, pp. 1724–1754, 2014.

- [8] E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?,” *Journal of the ACM (JACM)*, vol. 58, no. 3, pp. 1–37, 2011.
- [9] A. Elgabli, J. Park, A. S. Bedi, M. Bennis, and V. Aggarwal, “Gadmm: Fast and communication efficient framework for distributed machine learning.,” *Journal of Machine Learning Research*, vol. 21, no. 76, pp. 1–39, 2020.
- [10] S. Bi and Y.-J. A. Zhang, “An admm based method for computation rate maximization in wireless powered mobile-edge computing networks,” in *2018 IEEE International Conference on Communications (ICC)*, pp. 1–7, IEEE, 2018.