# OPTIMIZATION FOR MACHINE LEARNING
### Regularized, large-scale and distributed optimization
December 4, 2023

Today (last lecture!)
   Large-scale and decentralized optimization

Exam: Thursday December 14 (open book)

Projects: Due Friday January 19 AOE

# DUALITY AND DECENTRALIZED OPTIMIZATION

① A basic intro to (Lagrangian) duality

$\hookrightarrow$ Consider a linearly constrained optimization problem of the form

$$(P) \quad \underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) \quad \text{s.t. } Ax = b$$

where $f : \mathbb{R}^d \to \mathbb{R}$, $A \in \mathbb{R}^{m \times d}$ and $b \in \mathbb{R}^m$

We suppose that $(P)$ has a solution and we let

$$x^* \in \underset{x \in \mathbb{R}^d}{\arg\min} \{ f(x) \text{ s.t. } Ax = b \} \subseteq \mathbb{R}^d$$

$$f^* = \underset{x \in \mathbb{R}^d}{\min} \{ f(x) \text{ s.t. } Ax = b \} \in \mathbb{R}$$

$\hookrightarrow$ The Lagrangian of $(P)$ is the function

$$\mathcal{L} : \mathbb{R}^d \times \mathbb{R}^m \longrightarrow \mathbb{R}$$
$$(x, y) \longmapsto f(x) + y^T \overbrace{(Ax - b)}^{\in \mathbb{R}^m}$$

objective function

constraint $Ax - b = 0$

$\longrightarrow \mathcal{L}$ is a linear combination of the objective and the constraint functions

$y$ : Lagrange multipliers / dual variables

$\longrightarrow$ (P) is equivalent to

uncontrained problem (or with a complicated objective function) $\begin{cases} \text{minimize} & \max_{y \in \mathbb{R}^m} \mathcal{L}(x,y) \\ x \in \mathbb{R}^d \end{cases}$

↳ Solution to an optimization problem in $y$

$$\text{If } Ax - b \neq 0, \text{ then } \max_{y \in \mathbb{R}^m} \mathcal{L}(x,y) = \infty$$

$$\text{If } Ax - b = 0, \text{ then } \max_{y \in \mathbb{R}^m} \mathcal{L}(x,y) = f(x)$$

Def: The (Lagrangian) dual of (P) is the problem

(D) $\qquad \underset{y \in \mathbb{R}^m}{\text{maximize}} \quad \min_{x \in \mathbb{R}^d} \mathcal{L}(x,y)$

↑ dual variables

↓ Dual function

$y \mapsto \min_{x \in \mathbb{R}^d} \mathcal{L}(x,y)$

$= \min_{x \in \mathbb{R}^d} f(x) + y^T(Ax - b)$

• (D) is always a "convex problem"
   (= maximize the negative of a convex function)

$y \mapsto - \min_{x \in \mathbb{R}^d} \mathcal{L}(x,y)$ is a convex function

even when $f$ is not convex

- Let $q$ denote the dual function of the problem $\left(q(y) = \min\limits_{x \in \mathbb{R}^n} \mathcal{L}(x,y)\right)$ and let $q^* = \max\limits_{y \in \mathbb{R}^m} q(y)$.

Then, $$q^* \leq f^* \implies (D) \text{ gives an approximation of the optimal value of (P)}$$
$$(\text{"weak duality"})$$

$\longrightarrow$ In general, we cannot guarantee more than duality

$\longrightarrow$ But in our case, since the constraints are linear, we can guarantee that $q^* = f^*$ ("strong duality")

$$(\Rightarrow) \quad \exists \, y^* \in \mathbb{R}^m, \quad q(y^*) = q^* = f^* = f(x^*)$$

NB: In that case, we say that $(x^*, y^*)$ is a saddle point of $\mathcal{L}$

$$\forall \, x \in \{x \mid Ax = b\}, \quad \forall y \in \mathbb{R}^m,$$

$$\mathcal{L}(x^*, y) \leq \mathcal{L}(x^*, y^*) \leq \mathcal{L}(x, y^*)$$

$y^*$ maximizes $\mathcal{L}(x^*, \cdot)$

$q^* = f^*$

$x^*$ minimizes $\mathcal{L}(\cdot, y^*)$

$\longrightarrow$ With strong duality, we can build algorithms to solve (D) instead of (P) $\implies$ Dual algorithms

# (2) Dual algorithms

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \; f(x) \quad \text{s.t.} \quad Ax = b$$

## a) Dual ascent

(ascent $\Rightarrow$ $\underset{y \in \mathbb{R}^m}{\text{maximize}} \; q(y)$)

Algorithm ($x_0 \in \mathbb{R}^d$ not necessarily feasible, $y_0 \in \mathbb{R}^m$)

For $k = 0, 1, \text{---}$

- Compute $x_{k+1} \in \underset{x \in \mathbb{R}^d}{\text{argmin}} \; \mathcal{L}(x, y_k)$

- Set $y_{k+1} = y_k + \alpha_k \left( A x_{k+1} - b \right)$

  for some $\alpha_k > 0$

$x_{k+1} \in \underset{x \in \mathbb{R}^d}{\text{argmin}} \; \mathcal{L}(x, y_k)$: Solve an <u>unconstrained</u> optimization problem in $x$

$y_{k+1} = y_k + \alpha_k (A x_{k+1} - b)$: Subgradient step for the function $q : y \mapsto \underset{x \in \mathbb{R}^d}{\text{min}} \; \mathcal{L}(x, y)$ at $y_k$

$q$ is not differentiable in general but $-q$ is convex and it has subgradients at every point

Recall: $\underset{y \in \mathbb{R}^m}{\text{minimize}} \; -q(y) \qquad -q$ convex

Subgradient iteration $\qquad y_{k+1} = y_k - \alpha_k g_k, \quad g_k \in \partial(-q)(y_k)$

$\Rightarrow$ Analogous to (sub)gradient methods for the dual

$\Rightarrow$ Could be combined with stochastic (sub)gradient estimates, coordinate |ascent/|methods , acceleration, etc.
|descent|

## b) Augmented Lagrangian (aka method of multipliers)

$\longrightarrow$ In general, dual ascent will converge <u>slowly</u> to a solution (or a point with zero as a subgradient )

$\longrightarrow$ A dual ascent iteration is not uniquely defined and because of the ambiguity in choosing $x_{k+1}$ it can produce a sequence $\{x_k\}$ of infeasible points

$\Rightarrow$ Dual ascent only converges under some restrictions on the choice of $x_{k+1}$

$\Rightarrow$ Fix : use regularization to penalize infeasible points

Def: The augmented Lagrangian of (P) with parameter $\alpha > 0$ is defined as

$$\mathcal{L}_\alpha (x, y) = f(x) + y^T(Ax - b) + \frac{\alpha}{2} \|Ax - b\|^2$$

$$\rightarrow \quad \mathcal{L}_\lambda(x,y) = \mathcal{L}(x,y) + \frac{\lambda}{2} \|Ax-b\|^2$$

$\longrightarrow$ Lagrangian function of the regularized problem

$$\text{minimize}_{x \in \mathbb{R}^d} \quad f(x) + \frac{\lambda}{2}\|Ax-b\|^2$$

This problem is equivalent to (P) $\longrightarrow$

$$\text{s.t.} \qquad Ax=b$$

Remark: Many other augmented Lagrangian functions can be defined using other regularization terms (e.g. $\ell_1$)

Augmented Lagrangian algorithm ( $x_0 \in \mathbb{R}^d$, $y_0 \in \mathbb{R}^m$, $\lambda > 0$)

For $k=0,1,2,\ -$

• $x_{k+1} \in \text{argmin}_{x \in \mathbb{R}^d} \mathcal{L}_\lambda(x, y_k)$

• $y_{k+1} = y_k + \alpha_k \cdot (Ax_{k+1} - b)$

$\rightarrow$ Because of the regularization term, $x_{k+1}$ is more likely to be uniquely defined. (for instance, if $f$ is convex then $\mathcal{L}_\lambda(\cdot, y_k)$ is strongly convex $\Rightarrow$ unique minimum)

$\longrightarrow$ A popular choice for $\alpha_k$ is $\alpha_k = \lambda$, in which case we

have $\mathcal{L}_\lambda(x_{k+1}, y_{k+1}) = f(x_{k+1}) + y_{k+1}^T(Ax_{k+1}-b) + \frac{\lambda}{2}\|Ax_{k+1}-b\|^2$

$$y_{k+1}^T(Ax_{k+1}-b)$$
$$= [y_k + \lambda(Ax_{k+1} - b)]^T (Ax_{k+1}-b)$$

$$= f(x_{k+1}) + y_k^T(Ax_{k+1}-b)$$
$$+ \lambda(Ax_{k+1}-b)^T(Ax_{k+1}-b) \qquad + \frac{\lambda}{2}\|Ax_{k+1}-b\|^2$$

$$= y_k^T (A x_{k+1} - b)$$
$$+ \lambda (A x_{k+1} - b)^T (A x_{k+1} - b)$$
$$= y_k^T (A x_{k+1} - b)$$
$$+ \lambda \| A x_{k+1} - b \|^2$$

$$= \mathcal{L}_\lambda (x_{k+1}, y_k) + \lambda \| A x_{k+1} - b \|^2$$
$$> \mathcal{L}_\lambda (x_{k+1}, y_k)$$

Same magnitude than

$$\frac{\lambda}{2} \| A x_{k+1} - b \|^2$$

↳ Large step sizes are allowed in dual methods because they can be compensated at the next iteration by improvements towards feasibility ($Ax = b$)

↳ Other motivation for choosing $\alpha_k = \lambda$ : optimality conditions

$$x_{k+1} \in \operatorname*{argmin}_{x \in \mathbb{R}^d} \mathcal{L}(x_{k+1}, y_k)$$

when $f$ is differentiable, then we must have

Gradient with respect to the $x$ variables ⟶

$$\nabla_x \mathcal{L}_\lambda (x_{k+1}, y_k) = 0_{\mathbb{R}^d}$$

$$\nabla f(x_{k+1}) + A^T y_k + \lambda A^T (A x_{k+1} - b) = 0$$

$$\nabla f(x_{k+1}) + A^T \left[ y_k + \lambda (A x_{k+1} - b) \right] = 0$$

$$y_{k+1}$$

↳ If $x_{k+1}$ is feasible ($A x_{k+1} = b$), then $y_{k+1} = y_k$

and $x_{k+1} \in \operatorname*{argmin}_{x \in \mathbb{R}^d} \mathcal{L}_\lambda (x, y_k)$

$$\mathcal{L}_d(x_{n+1}, y_n) \le \mathcal{L}_d(x, y_n) \qquad \forall\, x \in \mathbb{R}^d$$

$$\Rightarrow \quad \mathcal{L}_d(x_{n+1}, y_n) \le \mathcal{L}_d(x, y_n) \qquad \forall\, x \in \mathbb{R}^d, \\ Ax = b$$

$$Ax = b \quad \Rightarrow \quad \mathcal{L}_d(x, y_k) = f(x) + y^T \underbrace{(Ax - b)}_{= 0} + \frac{d}{2} \|\underbrace{Ax - b}_{= 0}\|^2$$
$$= f(x)$$

Hence

$$\mathcal{L}_d(x_{n+1}, y_n) = f(x_{n+1}) \le \mathcal{L}_d(x, y_n) = f(x) \qquad \forall\, x \in \mathbb{R}^d, \\ Ax = b$$

$$\Rightarrow \quad x_{n+1} \in \underset{x \in \mathbb{R}^d}{\arg\min} \{ f(x) \text{ s.t. } Ax = b \}$$

(That argument also works for dual ascent)

---

③ Dual methods with decomposition

Decomposition: Use a specific problem structure to solve several small problems instead of a large one

For simplicity, we consider a decomposition in two "blocks"

$$(P_2) \quad \underset{\substack{u \in \mathbb{R}^{d_1} \\ v \in \mathbb{R}^{d_2}}}{\text{minimize}} \quad f_1(u) + f_2(v) \quad \text{s.t.} \quad A_1 u + A_2 v = b$$

$$f_1 : \mathbb{R}^{d_1} \longrightarrow \mathbb{R} \qquad A_1 \in \mathbb{R}^{m \times d_1}$$
$$f_2 : \mathbb{R}^{d_2} \longrightarrow \mathbb{R} \qquad A_2 \in \mathbb{R}^{m \times d_2} \qquad b \in \mathbb{R}^m$$

$\Longrightarrow$ Special case of $\quad$ minimize $f(x)$ s.t. $\quad Ax = b$
$\qquad\qquad\qquad\qquad\qquad x \in \mathbb{R}^d$

where $\quad f$ is partially separable

$$x = \begin{bmatrix} u \\ v \end{bmatrix} \begin{matrix} \updownarrow d_1 \\ \updownarrow d_2 \end{matrix} \qquad f(x) = f_1(u) + f_2(v)$$

$\Rightarrow$ No "coupling" between $u$ and
$v$ in the objective

$$A x = \begin{pmatrix} A_1 & A_2 \end{pmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = A_1 u + A_2 v$$

$\Longrightarrow$ linear coupling between $u$ and $v$
in the constraints

Q) How do we adapt the dual methods to take the
decomposition into account?

a) Dual ascent $\Longrightarrow$ Dual decomposition

$\hookrightarrow$ Variant of dual ascent that exploits the particular problem
structure

$\hookrightarrow$ The lagrangian function for $(P_2)$ is

$$\mathcal{L}(u, v, y) = f_1(u) + f_2(v) + y^T (A_1 u + A_2 v - b)$$

Dual ascent iteration for $(P_2)$:

$\bullet$ $(u_{k+1}, v_{k+1}) \in \underset{u, v}{\arg\min}\ \mathcal{L}(u, v, y_k)$

$\bullet$ $y_{k+1} = y_k + \alpha_k (A_1 u_{k+1} + A_2 v_{k+1} - b)$

Dual decomposition : uses the partial separability of $\mathcal{L}(u, v, y)$

Iteration $k$

$$u_{k+1} \in \underset{u \in \mathbb{R}^{d_1}}{\arg\min} \ \mathcal{L}(u, v_k, y_k)$$

$$v_{k+1} \in \underset{v \in \mathbb{R}^{d_2}}{\arg\min} \ \mathcal{L}(u_k, v, y_k)$$

$$y_{k+1} = y_k + \alpha_k (A_1 u_{k+1} + A_2 v_{k+1} - b)$$

* Smaller optimization problems to solve
* The two minimization problems can be solved in parallel

NB: this idea extends to multiple "blocks" of variables

$$\underset{\substack{x^{(1)} \in \mathbb{R}^{d_1} \\ \vdots \\ x^{(\ell)} \in \mathbb{R}^{d_\ell}}}{\text{minimize}} \quad \sum_{i=1}^{\ell} f_i(x^{(i)}) \quad \text{s.t.} \quad \sum_{i=1}^{\ell} A_i x^{(i)} = b$$

$d_1 + \ldots + d_\ell = d$

b) Alternating Direction Method of Multipliers (ADMM)

Idea : • Combine dual decomposition with augmented Lagrangian
• Similarly to block coordinate descent, perform updates over one block of variables at a time

Iteration $k$ of ADMM for $(P_2)$

( Augmented Lagrangian for $(P_2)$

$$\mathcal{L}_1(u, v, y) = f_1(u) + f_2(v) + y^\top (A_1 u + A_2 v - b) + \frac{1}{2} \| A_1 u + A_2 v - b \|^2$$

- $u_{k+1} \in \underset{u \in \mathbb{R}^{d_1}}{\text{argmin}} \; \mathcal{L}_\lambda (u, v_k, y_k)$

- $v_{k+1} \in \underset{v \in \mathbb{R}^{d_2}}{\text{argmin}} \; \mathcal{L}_\lambda (u_{k+1}, v, y_k)$ $\longrightarrow$ New value for $u$

- $y_{k+1} = y_k + \lambda (A_1 u_{k+1} + A_2 v_{k+1} - b)$

$\ominus$ The two minimization problems can no longer be parallelized

$\oplus$ We benefit from the update on $u$ when computing the new value for $v$

$\hookrightarrow$ Can show convergence results for ADMM, especially in the convex setting ($f_1, f_2$ convex). Those results are mainly asymptotic (no convergence rates), and have the form

$$\| A_1 u_{k+1} + A_2 v_{k+1} - b \| \underset{k \to \infty}{\longrightarrow} 0 \qquad \text{(} v \text{ towards feasibility)}$$

$$f_1(u_k) + f_2(v_k) \underset{k \to \infty}{\longrightarrow} f^\circ \text{ (optimum)}$$

$$y_k \underset{k \to \infty}{\longrightarrow} y^\circ \text{ (optimal dual variable)}$$

$\longrightarrow$ Since the 2010s, many variants on the ADMM have been proposed:

$\longrightarrow$ Stochastic ADMM

$\longrightarrow$ Proximal ADMM

$\longrightarrow$ Accelerated ADMM

$\longrightarrow$ Coordinate ADMM

$\longrightarrow$ ...

# ④ Application: Consensus and decentralized optimization ($\approx$ federated learning)

**Setup:**
$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \sum_{i=1}^{m} f_i(x), \qquad m >> 1 \text{ typically}$$

**Key assumption:** $\rightarrow$ The $f_i$'s (or the data used to compute the $f_i$'s) are distributed among $m$ entities or agents



$\rightarrow$ Every agent has a copy of $x$, denoted by $x^{(i)}$ and can improve $x^{(i)}$ with respect to its function $f_i$

$\rightarrow$ There is a global copy of $x$

$\rightarrow$ **Goal:** All agents must agree on the value of $x$ and it must be a solution of $\underset{x \in \mathbb{R}^d}{\text{minimize}} \sum_{i=1}^{m} f_i(x)$

## Consensus optimization problem

$$\underset{\substack{x \in \mathbb{R}^d \\ x^{(1)} \in \mathbb{R}^d \\ \vdots \\ x^{(m)} \in \mathbb{R}^d}}{\text{minimize}} \sum_{i=1}^{m} f_i\left(x^{(i)}\right) \qquad \text{s.t.} \quad \begin{array}{l} x - x^{(1)} = 0 \\ x - x^{(2)} = 0 \\ \vdots \\ x - x^{(m)} = 0 \end{array}$$

$\rightarrow$ We can apply ADMM on this problem with either 2 blocks $\left(x^{(1)}, \ldots, x^{(m)}\right)$ and $x$
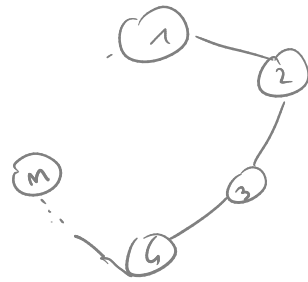
or $m+1$ blocks $\{x^{(1)}\}, \{x^{(2)}\}, \ldots, \{x^{(m)}\}, x$

# Decentralized optimization

$\hookrightarrow$ Same original problem, $m$ agents each with their own data ($f_i$) and their own copy of the variables $x^{(i)}$

$\hookrightarrow$ Agents organized in a network/graph $G = (V, E)$

$\quad V = \{1, \ldots, m\}$ vertices/agents

$\quad E \subseteq \{1, \ldots, m\} \times \{1, \ldots, m\}$ set of edges in the graphs

$\Rightarrow$ More general than the consensus problem: no central entity in general



# Decentralized optimization problem

$$\underset{\substack{x^{(1)} \in \mathbb{R}^d \\ \vdots \\ x^{(m)} \in \mathbb{R}^d}}{\text{minimize}} \quad \sum_{i=1}^{m} f_i(x^{(i)}) \quad \text{s.t.} \quad x^{(i_1)} - x^{(i_2)} = 0 \quad \forall (i_1, i_2) \in E$$

$\Rightarrow$ Partially separable: can apply dual decomposition or ADMM

$\Rightarrow$ Key for efficiency: reach consensus using

as little communication as possible between
agents