

First-Order Methods through Partial Linearization

Alp Yurtsever
Umeå University, Sweden

based on joint works with Suvrit Sra, Hooman Maskan, and others

May 2026
ELLIIT Focus Period
Lund, Sweden

DC Template

$$\min_{\mathbf{x} \in \mathbb{R}^m} \phi(\mathbf{x}) := f(\mathbf{x}) - h(\mathbf{x}), \quad (\text{DC})$$

- f and h are lower semi-continuous convex functions,
- we assume that the optimal value ϕ_* is finite.

DC Algorithm

The key idea of DCA (Pham Dinh & Souad, 1986) is to linearize the concave component $-h(\mathbf{x})$ to obtain a global convex upper bound:

$$\hat{\phi}(\mathbf{x}; \mathbf{y}) := f(\mathbf{x}) - h(\mathbf{y}) - \langle \nabla h(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle.$$

DCA then updates its estimation by solving convex subproblems

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^m} \hat{\phi}(\mathbf{x}; \mathbf{x}_k) \quad (1)$$

A solution to (1) always exists as ϕ_* is finite by assumption and

$$\min_{\mathbf{x} \in \mathbb{R}^m} \hat{\phi}(\mathbf{x}; \mathbf{x}_k) \geq \min_{\mathbf{x} \in \mathbb{R}^m} \phi(\mathbf{x}) = \phi_*.$$

DC Algorithm

- 1 DCA is a specific Majorize-Minorize algorithm.
- 2 If the functions are differentiable, DCA amounts to the implicit update rule:

$$\nabla f(\mathbf{x}^{k+1}) = \nabla h(\mathbf{x}^k).$$

- 3 DCA generates a monotonically decreasing sequence $\phi(\mathbf{x}^k)$.
- 4 Numerous algorithms can be viewed as special cases of DCA:
 - Sinkhorn's method for matrix scaling
 - EM algorithm for natural exponential family distributions
 - Proximal gradient method, Mirror descent
 - Proximal point method, Mirror prox

DC Critical Points

Definition

We measure convergence in terms of the following gap function:

$$\text{gap}_{\text{DC}}(\mathbf{y}) = \max_{\mathbf{x} \in \mathbb{R}^d} \min_{\mathbf{u} \in \partial h(\mathbf{y})} \{f(\mathbf{y}) - f(\mathbf{x}) - \langle \mathbf{u}, \mathbf{y} - \mathbf{x} \rangle\}.$$

Lemma

- $\text{gap}_{\text{DC}}(\mathbf{y}) \geq 0$ for all $\mathbf{y} \in \mathbb{R}^d$.
- $\text{gap}_{\text{DC}}(\mathbf{y}) = 0$ if and only if \mathbf{y} is a critical point satisfying

$$\partial f(\mathbf{y}) \cap \partial h(\mathbf{y}) \neq \emptyset.$$

If h is differentiable, then any such \mathbf{y} is a first-order stationary point.

Theoretical guarantees for DCA

Theorem

Let the sequence $\mathbf{x}^1, \dots, \mathbf{x}^K$ be generated by an inexact DCA, where each subproblem is solved approximately to satisfy

$$\hat{\phi}(\mathbf{x}^{k+1}, \mathbf{x}^k) - \min_{\mathbf{x} \in \mathbb{R}^d} \hat{\phi}(\mathbf{x}, \mathbf{x}^k) \leq \frac{\epsilon}{2}$$

for some target accuracy $\epsilon > 0$. Then,

$$\frac{1}{K} \sum_{k=1}^K \text{gap}_{\text{DC}}(\mathbf{x}^k) \leq \frac{\phi(\mathbf{x}^1) - \phi^*}{K} + \frac{\epsilon}{2}.$$

Implications:

- Iteration complexity of DCA is $\mathcal{O}(1/\epsilon)$.
- Complexity is independent of d (assuming $\phi(\mathbf{x}^1) - \phi^*$ is constant).

Proof: Theoretical guarantees for DCA

Denoting $\mathbf{u}^k \in \partial h(\mathbf{x}^k)$, by the convexity of h ,

$$\phi(\mathbf{x}^{k+1}) = f(\mathbf{x}^{k+1}) - h(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^{k+1}) - h(\mathbf{x}^k) - \langle \mathbf{u}^k, \mathbf{x}^{k+1} - \mathbf{x}^k \rangle.$$

Then, by definition of \mathbf{x}^{k+1} :

$$\phi(\mathbf{x}^{k+1}) \leq f(\mathbf{x}) - h(\mathbf{x}^k) - \langle \nabla h(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \epsilon/2, \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

Add and subtract $f(\mathbf{x}^k)$ on the right side:

$$\begin{aligned} \phi(\mathbf{x}^{k+1}) &\leq f(\mathbf{x}^k) - f(\mathbf{x}^k) + f(\mathbf{x}) - h(\mathbf{x}^k) - \langle \nabla h(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \epsilon/2 \\ &\leq \phi(\mathbf{x}^k) - \text{gap}_{\text{DC}}(\mathbf{x}^k) + \epsilon/2. \end{aligned}$$

We complete the proof by averaging over $k = 1, \dots, K$.

⚠ Convexity of f is not used!

$$\min_{\mathbf{x} \in \mathbb{R}^m} F(\mathbf{x}).$$

Add and subtract a strongly convex functions:

$$\min_{\mathbf{x} \in \mathbb{R}^m} \underbrace{F(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x}\|^2}_{f(\mathbf{x})} - \underbrace{\frac{1}{2\lambda} \|\mathbf{x}\|^2}_{h(\mathbf{x})}.$$

DCA update is

$$\begin{aligned} \mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^m} \left\{ F(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x}\|^2 - \frac{1}{\lambda} \langle \mathbf{x}_k, \mathbf{x} - \mathbf{x}_k \rangle \right\} \\ &= \operatorname{prox}_{\lambda F}(\mathbf{x}_k) \end{aligned}$$

DCA convergence theorem implies

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \max_{\mathbf{x} \in \mathbb{R}^m} \left\{ F(\mathbf{x}_k) - F(\mathbf{x}) - \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{x}_k\|^2 \right\} &\leq \frac{\phi(\mathbf{x}_1) - \phi_\star}{K} \\ \implies \frac{1}{K} \sum_{k=1}^K \frac{1}{2(L + \lambda^{-1})} \|\mathbf{x}_k - \operatorname{prox}_{\lambda F}(\mathbf{x}_k)\|^2 &\leq \frac{\phi(\mathbf{x}_1) - \phi_\star}{K} \end{aligned}$$

PGM via DCA

$$\min_{\mathbf{x} \in \mathbb{R}^m} F(\mathbf{x}) + G(\mathbf{x})$$

Add and subtract a strongly convex functions:

$$\min_{\mathbf{x} \in \mathbb{R}^m} \underbrace{G(\mathbf{x}) + \frac{L}{2}\|\mathbf{x}\|^2}_{f(\mathbf{x})} - \underbrace{\left(\frac{L}{2}\|\mathbf{x}\|^2 - F(\mathbf{x})\right)}_{h(\mathbf{x})}$$

DCA update is

$$\begin{aligned}\mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^m} \left\{ G(\mathbf{x}) + \frac{L}{2}\|\mathbf{x}\|^2 - \langle L\mathbf{x}_k - \nabla F(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle \right\} \\ &= \operatorname{prox}_{\frac{1}{L}G}(\mathbf{x}_k - \frac{1}{L}\nabla F(\mathbf{x}_k))\end{aligned}$$

DCA convergence theorem implies

$$\begin{aligned}\frac{1}{K} \sum_{k=1}^K \max_{\mathbf{x} \in \mathbb{R}^m} \left\{ \langle \nabla F(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x} \rangle + G(\mathbf{x}_k) - G(\mathbf{x}) - \frac{L}{2}\|\mathbf{x} - \mathbf{x}_k\|^2 \right\} &\leq \frac{\phi(\mathbf{x}_1) - \phi_\star}{K} \\ \implies \frac{1}{K} \sum_{k=1}^K \frac{L}{2} \|\mathbf{x}_k - \operatorname{prox}_{\frac{1}{L}G}(\mathbf{x}_k - \frac{1}{L}\nabla F(\mathbf{x}_k))\|^2 &\leq \frac{\phi(\mathbf{x}_1) - \phi_\star}{K}\end{aligned}$$

EM via DCA

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^m} \mathcal{L}(\boldsymbol{\theta}) := - \sum_{x \in \mathcal{X}} \log P(x|\boldsymbol{\theta})$$

P is defined through hidden variables: $P(x|\boldsymbol{\theta}) = \sum_y P(x, y|\boldsymbol{\theta})$.

Consider the natural exponential family of distributions, defined by

$$P(x, y|\boldsymbol{\theta}) = \frac{\varphi(x, y)e^{\langle \boldsymbol{\theta}, T(x, y) \rangle}}{\sum_{x, y} \varphi(x, y)e^{\langle \boldsymbol{\theta}, T(x, y) \rangle}},$$

EM can be viewed as DCA applied to

$$f(\boldsymbol{\theta}) = |\mathcal{X}| \log \left(\sum_{x, y} \varphi(x, y)e^{\langle \boldsymbol{\theta}, T(x, y) \rangle} \right)$$
$$h(\boldsymbol{\theta}) = \sum_{x \in \mathcal{X}} \log \left(\sum_y \varphi(x, y)e^{\langle \boldsymbol{\theta}, T(x, y) \rangle} \right)$$

These terms are smooth, owing to log-sum-exp structure.

DCA convergence theorem implies

$$\frac{1}{K} \sum_{k=1}^K \max_{\boldsymbol{\theta} \in \mathbb{R}^m} \{f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}) - \langle \nabla h(\boldsymbol{\theta}_k), \boldsymbol{\theta}_k - \boldsymbol{\theta} \rangle\} \leq \frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}_*}{K}$$

$$\implies \frac{1}{K} \sum_{k=1}^K \max_{\boldsymbol{\theta} \in \mathbb{R}^m} \left\{ \underbrace{\langle \nabla f(\boldsymbol{\theta}_k) - \nabla h(\boldsymbol{\theta}_k), \boldsymbol{\theta}_k - \boldsymbol{\theta} \rangle}_{\nabla \mathcal{L}(\boldsymbol{\theta}_k)} - \frac{L}{2} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}\|^2 \right\} \leq \frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}_*}{K}$$

$$\iff \frac{1}{K} \sum_{k=1}^K \frac{1}{2L} \|\nabla \mathcal{L}(\boldsymbol{\theta}_k)\|^2 \leq \frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}_*}{K}$$

▲ This rate is consistent with the existing guarantees, as provided by [Kumar and Schmidh \(2017\)](#) and [Kunstner et al. \(2020\)](#).

DCA is FW in disguise

DCA is FW

Proposition

DCA is equivalent to the FW method applied to the following epigraph reformulation of (DC):

$$\min_{\mathbf{x}, t} \quad t - h(\mathbf{x}), \quad f(\mathbf{x}) \leq t.$$

$$\underbrace{\min_{\mathbf{x}, t}}_{\omega} \quad \underbrace{t - h(\mathbf{x})}_{\psi(\omega)}, \quad \underbrace{f(\mathbf{x}) \leq t}_{\omega \in \mathcal{M}}.$$

The linear minimization step of FW amounts to

$$(\mathbf{x}_k^*, t_k^*) \in \operatorname{argmin}_{\mathbf{x}, t} t - \langle \nabla h(\mathbf{x}_k), \mathbf{x} \rangle \quad \text{s.t.} \quad f(\mathbf{x}) \leq t. \quad (2)$$

Eliminate t in (2) and get

$$\mathbf{x}_k^* \in \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) - \langle \nabla h(\mathbf{x}_k), \mathbf{x} \rangle.$$

FW update becomes

$$\mathbf{x}_{k+1} = (1 - \eta_k) \mathbf{x}_k + \eta_k \mathbf{x}_k^* = \mathbf{x}_k^*$$

Curvature constant

The convergence of FW crucially depends on the *curvature constant* of ψ over \mathcal{M} , defined as follows:

$$C_\psi := \sup_{\substack{\omega, \hat{\omega} \in \mathcal{M}; \eta \in [0, 1] \\ \bar{\omega} = (1-\eta)\omega + \eta\hat{\omega}}} \frac{2}{\eta^2} \left(\psi(\bar{\omega}) - \psi(\omega) - \langle \nabla \psi(\omega), \bar{\omega} - \omega \rangle \right).$$

Equivalently, for all $\omega, \hat{\omega} \in \mathcal{M}$; $\eta \in [0, 1]$, and $\bar{\omega} = (1 - \eta)\omega + \eta\hat{\omega}$;

$$\psi(\bar{\omega}) \leq \psi(\omega) + \langle \nabla \psi(\omega), \bar{\omega} - \omega \rangle + \frac{1}{2} \eta^2 C_\psi$$

Special cases:

- If ψ is L -smooth and \mathcal{M} has a bounded diameter D , then $C_\psi \leq LD^2$.
- If ψ is concave, then $C_\psi = 0$.

Guarantees of FW

Theorem (Jaggi, 2013, Theorem 1)

Suppose ψ is **convex** and C_ψ is bounded. Then, the sequence ω_k generated by FW converges to a solution in objective value:

$$\psi(\omega_k) - \psi_\star \leq \frac{2C_\psi}{k+1}.$$

Theorem (Lacoste-Julien, 2016, Theorem 1)

Suppose C_ψ is bounded but ψ is **non-convex**. Then,

$$\frac{1}{k} \sum_{\tau=1}^k \max_{\omega \in \mathcal{M}} \langle \nabla \psi(\omega_\tau), \omega_\tau - \omega \rangle \leq \frac{\max\{C_\psi, 2(\psi(\omega_1) - \psi_\star)\}}{\sqrt{k}}$$

Theorem (Journée et al., 2010, Theorem 1)

Suppose ψ is **concave**. Then,

$$\frac{1}{k} \sum_{\tau=1}^k \max_{\omega \in \mathcal{M}} \langle \nabla \psi(\omega_\tau), \omega_\tau - \omega \rangle \leq \frac{\psi(\omega_1) - \psi_\star}{k}.$$

Beyond DC decompositions

Classical DCA focuses on a convex – convex decomposition.

FW viewpoint suggests a broader range of decompositions.

$$\min_{\mathbf{x}, t} \quad t - h(\mathbf{x}), \quad f(\mathbf{x}) \leq t.$$

If h is non-convex:

- full step $\eta = 1$ is no longer valid in general,
- the rate reduces to $\mathcal{O}(1/\sqrt{k})$.

The extended algorithm is

$$\mathbf{s}_k \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^m} \hat{\phi}(\mathbf{x}; \mathbf{x}_k)$$
$$\mathbf{x}_{k+1} = (1 - \eta_k)\mathbf{x}_k + \eta_k \mathbf{s}_k$$

Subgradient Method via DCA

$$\min_{\mathbf{x} \in \mathbb{R}^m} F(\mathbf{x})$$

Add and subtract a strongly convex function:

$$\min_{\mathbf{x} \in \mathbb{R}^m} \underbrace{\frac{1}{2\alpha} \|\mathbf{x}\|^2}_{f(\mathbf{x})} - \underbrace{\left(\frac{1}{2\alpha} \|\mathbf{x}\|^2 - F(\mathbf{x}) \right)}_{h(\mathbf{x})}$$

h can be nonconvex, since F is convex but nonsmooth.

Let $\mathbf{u}_k \in \partial F(\mathbf{x}_k)$. DCA update is

$$\mathbf{s}_k = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^m} \left\{ \frac{1}{2\alpha} \|\mathbf{x}\|^2 - \langle \frac{1}{\alpha} \mathbf{x}_k - \mathbf{u}_k, \mathbf{x} - \mathbf{x}_k \rangle \right\} = \mathbf{x}_k - \alpha \mathbf{u}_k$$

$$\mathbf{x}_{k+1} = (1 - \eta_k) \mathbf{x}_k + \eta_k \mathbf{s}_k = \mathbf{x}_k - \alpha \eta_k \mathbf{u}_k$$

A common step-size for nonconvex FW is $\eta_k = 1/\sqrt{k}$.

Dynamic Decompositions

Classical DCA focuses on a fixed decomposition

$$\phi(\mathbf{x}) := f(\mathbf{x}) - h(\mathbf{x}).$$

The analysis allows for a dynamic decomposition

$$\phi(\mathbf{x}) := f_k(\mathbf{x}) - h_k(\mathbf{x}).$$

⚠ $\text{gap}_{\text{DC}}(\mathbf{x}_k) := \text{gap}_{\text{DC}}^{f_k, h_k}(\mathbf{x}_k)$ depends on the decompositions.

We get the guarantees in terms of $\frac{1}{K} \sum_{k=1}^K \text{gap}_{\text{DC}}^{f_k, h_k}(\mathbf{x}_k)$.

Non-Euclidean GM via DCA

$$\min_{\mathbf{x} \in \mathbb{R}^m} F(\mathbf{x})$$

Add and subtract squared norms centered at \mathbf{x}_k :

$$\min_{\mathbf{x} \in \mathbb{R}^m} \underbrace{\frac{1}{2\alpha} \|\mathbf{x} - \mathbf{x}_k\|^2}_{f_k(\mathbf{x})} - \underbrace{\left(\frac{1}{2\alpha} \|\mathbf{x} - \mathbf{x}_k\|^2 - F(\mathbf{x}) \right)}_{h_k(\mathbf{x})}$$

DCA update is

$$\mathbf{s}_k = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^m} \left\{ \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{x}_k\|^2 - \frac{1}{\alpha} \langle \nabla F(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle \right\}$$

$$\mathbf{x}_{k+1} = (1 - \eta_k) \mathbf{x}_k + \eta_k \mathbf{s}_k.$$

Revisiting Frank-Wolfe for Structured Nonconvex Problems

FW for DC functions

$$\min_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x}) := f(\mathbf{x}) - h(\mathbf{x}),$$

- $\mathcal{D} \subset \mathbb{R}^m$ is a convex and compact set with diameter D ,
- f and h are lower semi-continuous convex functions,
- f is L_f -smooth.

DCFW Algorithm

Idea: Apply DCA, and solve the convex subproblems with FW.

$$\mathbf{x}_{t+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}) - \langle \mathbf{u}_t, \mathbf{x} \rangle, \quad \mathbf{u}_t \in \partial h(\mathbf{x}_t).$$

- 1: **Input:** initial point $\mathbf{x}_1 \in \mathcal{D}$, target accuracy $\epsilon > 0$
- 2: **for** $t = 1, 2, \dots$ **do**
- 3: Choose $\mathbf{u}_t \in \partial h(\mathbf{x}_t)$
- 4: Initialize $\mathbf{X}_{t,1} = \mathbf{x}_t$
- 5: **for** $k = 1, 2, \dots$ **do**
- 6: $\mathbf{S}_{t,k} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} \langle \nabla f(\mathbf{X}_{t,k}) - \mathbf{u}_t, \mathbf{x} \rangle$
- 7: $\mathbf{D}_{t,k} = \mathbf{S}_{t,k} - \mathbf{X}_{t,k}$
- 8: **if** $-\langle \nabla f(\mathbf{X}_{t,k}) - \mathbf{u}_t, \mathbf{D}_{t,k} \rangle \leq \epsilon/2$ **then**
- 9: Set $\mathbf{x}_{t+1} = \mathbf{X}_{t,k}$ and **break**
- 10: **end if**
- 11: $\mathbf{X}_{t,k+1} = \mathbf{X}_{t,k} + \eta_{t,k} \cdot \mathbf{D}_{t,k}$ // $\eta_{t,k} = 2/(k+1)$ or line-search
- 12: **end for**
- 13: **end for**

DCFW guarantees

Theorem

DCFW generates a sequence of solutions that satisfy

$$\min_{0 \leq \tau \leq t} \text{gap}_{\text{DC}}(\mathbf{x}_\tau) \leq \epsilon$$

within $\mathcal{O}(1/\epsilon)$ outer iterations, with a total of $\mathcal{O}(1/\epsilon^2)$ inner iterations.

**Intuition:* DCA converges at $\mathcal{O}(1/t)$ rate. Each DCA step solves a smooth convex problem via FW with $\mathcal{O}(1/k)$ rate.

Oracle Complexity to reach ϵ -accuracy:

- $\mathcal{O}(1/\epsilon^2)$ linear minimization oracles,
- $\mathcal{O}(1/\epsilon^2)$ gradient calls for f ,
- $\mathcal{O}(1/\epsilon)$ (sub)gradient calls for h .

Conditional gradient sliding

DC decomposition is not unique — choosing a suitable decomposition can reduce oracle complexity.

Example: Since f is L_f -smooth, $\frac{L_f}{2}\|\mathbf{x}\|^2 - f(\mathbf{x})$ is convex. Consider DCFW on the following reformulation of the problem:

$$\min_{\mathbf{x} \in \mathcal{D}} \underbrace{\frac{L_f}{2}\|\mathbf{x}\|^2}_{f'(\mathbf{x})} - \underbrace{\left(h(\mathbf{x}) + \frac{L_f}{2}\|\mathbf{x}\|^2 - f(\mathbf{x})\right)}_{h'(\mathbf{x})}$$

Benefit: Gradients of f are now accessed less frequently.

Oracle Complexity to reach ϵ -accuracy:

- $\mathcal{O}(1/\epsilon^2)$ linear minimization oracle calls,
- $\mathcal{O}(1/\epsilon)$ gradient calls for f ,
- $\mathcal{O}(1/\epsilon)$ (sub)gradient calls for h .

When the domain is strongly convex

The inner FW loop can be faster if both the surrogate and the feasible region have curvature.

- Strongly convex sets contain a small ball around every chord.
- Examples: ℓ_p -norm balls and Schatten- p balls for $1 < p \leq 2$.

Theorem

If f is strongly convex and \mathcal{D} is strongly convex, then DCFW reaches

$$\min_{0 \leq \tau \leq t} \text{gap}_{\text{DC}}(\mathbf{x}_\tau) \leq \epsilon$$

within $\mathcal{O}(1/\epsilon)$ outer iterations and $\mathcal{O}(1/\epsilon^{3/2})$ linear minimization oracle calls.

Reason: FW has an $\mathcal{O}(1/k^2)$ rate for smooth strongly convex objectives over strongly convex sets.

Case study: QAP

Relax-and-round for the quadratic assignment problem:

$$\min_{\mathbf{X}} \langle \mathbf{A}^\top \mathbf{X}, \mathbf{X} \mathbf{B} \rangle \quad \text{s. t.} \quad \mathbf{X} \in [0, 1]^{n \times n}, \mathbf{X} \mathbf{1} = \mathbf{X}^\top \mathbf{1} = \mathbf{1}.$$

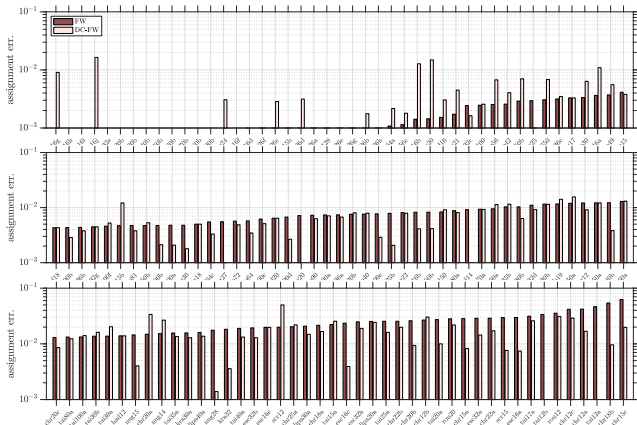
The feasible set is the Birkhoff polytope; the linear minimization oracle is a linear assignment problem.

DCFW uses the decomposition

$$f(\mathbf{X}) = \frac{1}{4} \|\mathbf{A}^\top \mathbf{X} + \mathbf{X} \mathbf{B}\|_{\mathbb{F}}^2, \quad h(\mathbf{X}) = \frac{1}{4} \|\mathbf{A}^\top \mathbf{X} - \mathbf{X} \mathbf{B}\|_{\mathbb{F}}^2.$$

QAPLIB instances	DCFW better	FW better	tie
134 total	73	43	18

QAP numerical result



Assignment error after rounding. Instances are ordered from best to worst performance of FW. DCFW also achieved a lower average assignment error: 0.0085 vs. 0.0112.

Block Coordinate DCA

Motivation: QBO

Quadratic Binary Optimization:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \lambda \|\mathbf{x}\|_1 + \lambda d \quad \text{subject to} \quad \mathbf{x} \in \{-1, 1\}^d$$

Convex hull relaxation:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \lambda \|\mathbf{x}\|_1 + \lambda d \quad \text{subject to} \quad \|\mathbf{x}\|_\infty \leq 1$$

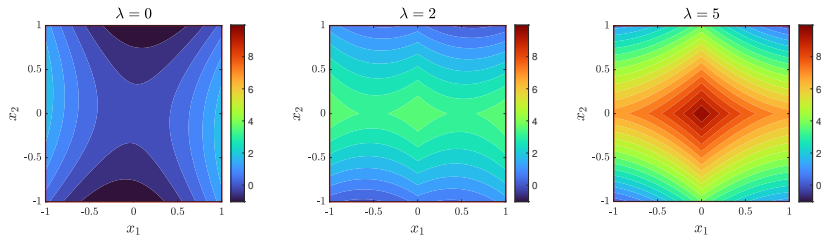


Figure: Illustration of level curves with increasing values of λ .

Negative ℓ_1 regularization promotes solutions closer to the binary set. We use $\mathbf{Q} = \begin{bmatrix} 2 & -0.25 \\ -0.25 & -1 \end{bmatrix}$.

Motivation: Challenge

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \lambda \|\mathbf{x}\|_1 \quad \text{subject to} \quad \|\mathbf{x}\|_\infty \leq 1$$

How to decompose into f and h ?

Approach I:

Split $\mathbf{Q} = \mathbf{Q}_P + \mathbf{Q}_N$ (positive and negative semidefinite parts) and

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}_P \mathbf{x} + \iota_{\{\|\mathbf{x}\|_\infty \leq 1\}}(\mathbf{x}), \quad h(\mathbf{x}) = -\mathbf{x}^\top \mathbf{Q}_N \mathbf{x} + \lambda \|\mathbf{x}\|_1.$$

Challenge:

- Requires eigendecomposition of \mathbf{Q} ; expensive for large d .
- \mathbf{Q}_P and \mathbf{Q}_N may lose useful structural properties (e.g., sparsity).

Approach II:

Consider nonconvex f by choosing

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \iota_{\{\|\mathbf{x}\|_\infty \leq 1\}}(\mathbf{x}), \quad h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1.$$

Challenge:

Motivation

Coordinate descent idea: Minimize over one coordinate (or block coordinate) at a time, fixing all other coordinates.

Notation:

- $x_i \in \mathbb{R}^{d_i}$ represents the i^{th} coordinate block of x .
- $\mathbf{x}_i \in \mathbb{R}^d$ is an extension of x_i (other coordinates padded with zeros).
- $\bar{\mathbf{x}}_i \in \mathbb{R}^d$ is the complement of \mathbf{x}_i , containing zeros in the i^{th} block and ensuring $\mathbf{x}_i + \bar{\mathbf{x}}_i = \mathbf{x}$.

Subproblem becomes:

$$x_{i_k}^{k+1} \leftarrow \operatorname{argmin}_{x_{i_k} \in [-1,1]} (\bar{\mathbf{x}}_{i_k}^k + \mathbf{x}_{i_k})^\top \mathbf{Q}(\bar{\mathbf{x}}_{i_k}^k + \mathbf{x}_{i_k}) - \lambda \langle \operatorname{sign}(\mathbf{x}^k), \bar{\mathbf{x}}_{i_k}^k + \mathbf{x}_{i_k} \rangle$$

This reduces to a simple quadratic subproblem:

$$x_i^{k+1} \leftarrow \operatorname{argmin}_{x_i \in [-1,1]} ax_i^2 + bx_i + c,$$

- $a = q_{ii}$ (q_{ii} is the i th diagonal entry of \mathbf{Q})
- $b = 2\mathbf{q}_i^\top \bar{\mathbf{x}}_i^k - \lambda \operatorname{sign}(x_i^k)$ (\mathbf{q}_i is the i th column of \mathbf{Q})

Problem template

$$\min_{\mathbf{x} \in \mathbb{R}^d} \phi(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) - h(\mathbf{x}), \quad (\text{DC})$$

- f, g and h are lower semi-continuous functions. h is convex.
- f is L -smooth, g can be nonsmooth (can take $+\infty$).
- g is (block) coordinate-wise separable: $g(\mathbf{x}) = \sum_{i=1}^n g_i(x_i)$.
- We assume that the optimal value ϕ^* is finite.

Block coordinate DCA

- 1 Start from a feasible initial point $\mathbf{x}^1 \in \mathbb{R}^d$.
- 2 For $k = 1, \dots, K$, update the estimate as follows:
 - a Choose a block i_k uniformly at random from $\{1, \dots, n\}$.
 - b Keep all other coordinates fixed, except i_k , i.e., $\bar{\mathbf{x}}_{i_k}^{k+1} = \bar{\mathbf{x}}_{i_k}^k$.
 - c Update the i_k^{th} block by minimizing the surrogate objective along these coordinates:

$$\mathbf{x}_{i_k}^{k+1} = \operatorname{argmin}_{\mathbf{x}_{i_k} \in \mathbb{R}^{d_i}} \hat{\phi}(\mathbf{x}, \mathbf{x}^k) \quad \text{subj.to} \quad \mathbf{x} = \mathbf{x}_{i_k} + \bar{\mathbf{x}}_{i_k}^k$$

Let us denote the surrogate objective restricted to i^{th} block by

$$\begin{aligned} \hat{\phi}_i(\mathbf{x}_i, \mathbf{x}^k) &:= \hat{\phi}(\mathbf{x}_i + \bar{\mathbf{x}}_i^k, \mathbf{x}^k) \\ &= f(\mathbf{x}_i + \bar{\mathbf{x}}_i^k) + g_i(\mathbf{x}_i) - g_i(\bar{\mathbf{x}}_i^k) + g(\mathbf{x}^k) - h(\mathbf{x}^k) - \langle \mathbf{u}_i^k, \mathbf{x}_i - \bar{\mathbf{x}}_i^k \rangle. \end{aligned}$$

where \mathbf{u}_i^k is the i^{th} block of a subgradient $\mathbf{u}^k \in \partial h(\mathbf{x}^k)$.

(Modified) Gap function

Definition (DC Gap)

We measure convergence in terms of the following gap function:

$$\text{gap}_{\text{DC}}(\mathbf{y}) = \max_{\mathbf{x} \in \mathbb{R}^d} \min_{\mathbf{u} \in \partial h(\mathbf{y})} \{f(\mathbf{y}) - f(\mathbf{x}) + g(\mathbf{y}) - g(\mathbf{x}) - \langle \mathbf{u}, \mathbf{y} - \mathbf{x} \rangle\}.$$

Definition (BDC Gap)

Consider the following gap function:

$$\text{gap}_{\text{BDC}}^L(\mathbf{y}) = \max_{\mathbf{x} \in \mathbb{R}^d} \min_{\mathbf{u} \in \partial h(\mathbf{y})} \{ \langle \nabla f(\mathbf{y}) - \mathbf{u}, \mathbf{y} - \mathbf{x} \rangle + g(\mathbf{y}) - g(\mathbf{x}) - \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \}$$

Lemma

- $\text{gap}_{\text{BDC}}^L(\mathbf{y}) \geq 0$ for all $\mathbf{y} \in \mathbb{R}^d$,
- $\text{gap}_{\text{BDC}}^L(\mathbf{y}) = 0$ if and only if \mathbf{y} is a critical point:

$$(\nabla f(\mathbf{y}) + \partial g(\mathbf{y})) \cap \partial h(\mathbf{y}) \neq \emptyset.$$

If h is differentiable, then any such \mathbf{y} is a first-order stationary point.

Convergence guarantees of BDCA

Theorem

Suppose $\mathbf{x}^1, \dots, \mathbf{x}^K$ is a sequence generated by BDCA, where each subproblem is solved approximately to satisfy

$$\hat{\phi}_{i_k}(\mathbf{x}_{i_k}^{k+1}, \mathbf{x}^k) - \min_{x_{i_k} \in \mathbb{R}^{d_i}} \hat{\phi}_{i_k}(x_{i_k}, \mathbf{x}^k) \leq \frac{\epsilon}{2}$$

for some target accuracy $\epsilon > 0$. Then,

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \left[\text{gap}_{\text{BDC}}^L(\mathbf{x}^k) \right] \leq \frac{n}{K} (\phi(\mathbf{x}^1) - \phi^*) + \frac{\epsilon}{2}.$$

Numerical experiment (methods)

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \lambda \|\mathbf{x}\|_1 \quad \text{subject to} \quad \|\mathbf{x}\|_\infty \leq 1$$

We consider 3 methods:

(Method 1) BDCA with nonconvex quadratic f :

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}, \quad g(\mathbf{x}) = \iota_{\{\|\mathbf{x}\|_\infty \leq 1\}}(\mathbf{x}), \quad h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$$

(Method 2) DCA with eigensplitting:

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}_P \mathbf{x} + \iota_{\{\|\mathbf{x}\|_\infty \leq 1\}}(\mathbf{x}), \quad h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1 - \mathbf{x}^\top \mathbf{Q}_N \mathbf{x}$$

(Method 3) BDCA with quadratic majorization:

$$f(\mathbf{x}) = \frac{L}{2} \|\mathbf{x}\|^2, \quad g(\mathbf{x}) = \iota_{\{\|\mathbf{x}\|_\infty \leq 1\}}(\mathbf{x}), \quad h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1 + \frac{L}{2} \|\mathbf{x}\|^2 - \mathbf{x}^\top \mathbf{Q} \mathbf{x}$$

where $L = 2\|\mathbf{Q}\|$.

Numerical experiment (results)

Experimental setup:

- $\lambda = \|\mathbf{Q}\|_F / \sqrt{d}$
- Initialization: same Gaussian *iid* vector for all methods
- Stop when gap reaches 10^{-6}
- GSet graph instances (67 datasets, d ranging from 800 to 10'000)

Table: Rounded objective values and running times for different methods on the largest 5 GSet instances.

Dataset	Size (d)	OBJ			TIME (s)		
		M1	M2	M3	M1	M2	M3
G63	7000	-107010.03	-61763	-107010.03	117.79	4218.65	354.19
G64	7000	-49222.03	-48250.03	-48670.03	43.75	1535.71	365.12
G65	8000	-32372	-32156	-32532	17.69	3468.09	32.90
G66	9000	-36720	-36360	-36704	25.39	1011.99	55.51
G67	10000	-40428	-40052	-40402	33.96	6298.46	63.50

- ⚠** Out of 67 problems, M1 finds the best solution in 56 instances, M2 in 22 instances, and M3 in 32 instances.

References

- Yurtsever, A. & Sra, S. (2022). CCCP is Frank-Wolfe in disguise.
- Maskan, H. Halvachi, P., Sra, S., & Yurtsever, A. (2024). Randomized block coordinate DC algorithm.
- Maskan, H., Hou, Y., Sra, S., & Yurtsever, A. (2025). Revisiting Frank-Wolfe for structured nonconvex optimization.
- Pham Dinh, T., & Souad, E.B. (1986). Algorithms for solving a class of nonconvex optimization problems: methods of subgradient.
- Frank, M., & Wolfe, P. (1956). An algorithm for quadratic programming.
- Jaggi, M. (2013). Revisiting Frank-Wolfe: Projection-free sparse convex optimization.
- Lacoste-Julien, S. (2016). Convergence rate of Frank-Wolfe for non-convex objectives.
- Journée, M., Nesterov, Y., Richtárik, P., & Sepulchre, R. (2010). Generalized Power Method for Sparse Principal Component Analysis.
- Kumar, R., & Schmidt, M. (2017) Convergence rate of expectation maximization.
- Kunstner, F., Kumar, R., & Schmidt, M. (2020) Homeomorphic-invariance of EM: Non-asymptotic convergence in KL divergence for exponential families via mirror descent.