

# Optimization Algorithms

Daniel P. Palomar (2025). *Portfolio Optimization: Theory and Application*.  
Cambridge University Press.

[portfoliooptimizationbook.com](https://portfoliooptimizationbook.com)

Latest update: 2025-09-15

# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary

# Executive Summary

- Over the past century, convex optimization has evolved through key algorithmic breakthroughs, beginning with Dantzig's 1947 simplex method for linear programming.
- Despite exponential worst-case complexity, the simplex method became widely adopted and laid the foundation for modern optimization.
- Karmarkar's 1984 interior-point method revolutionized linear programming by achieving polynomial time complexity, spurring extensive research into interior-point approaches.
- In 1994, Nesterov and Nemirovskii's theory of self-concordant functions enabled log-barrier algorithms to tackle broader convex problems including semidefinite programming and second-order cone programming.
- Specialized techniques such as block-coordinate descent, majorization-minimization, and successive convex approximation have emerged to create customized algorithms for specific problem structures.
- These specialized methods often provide enhanced complexity guarantees and improved convergence rates compared to general-purpose solvers.
- These slides explore a comprehensive range of practical algorithms developed through this rich algorithmic evolution (Palomar 2025, Appendix B).

# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary

- A solver, or optimizer, is an engine designed to solve specific types of mathematical problems.
- Available in various programming languages: R, Python, Matlab, Julia, Rust, C, C++.
- Each solver typically handles specific problem categories: LP, QP, QCQP, SOCP, SDP.

## **GLPK (GNU Linear Programming Kit):**

- Intended for large-scale LP including mixed-integer variables.
- Written in C.

## **Quadprog:**

- Popular open-source QP solver.
- Originally written in Fortran by Berwin Turlach in the late 1980s.
- Accessible from most programming languages.

## **MOSEK:**

- Proprietary solver for LP, QP, SOCP, SDP including mixed-integer variables.
- Established in 1997 by Erling Andersen.
- Specialized in large-scale problems; very fast, robust, and reliable.
- Free license available for academia.

## **SeDuMi:**

- Open-source solver for LP, QP, SOCP, SDP.
- Originally developed by Sturm in 1999 for Matlab.

## **SDPT3:**

- Open-source solver for LP, QP, SOCP, SDP.
- Originally developed in 1999 for Matlab.

## **Gurobi:**

- Proprietary solver for LP, QP, and SOCP including mixed-integer variables.
- Free license available for academia.

## **Embedded Conic Solver (ECOS):**

- SOCP solver originally written in C.

## **CPLEX:**

- Proprietary solver for LP and QP, also handles mixed-integer variables.
- Free license available for academia.

# Complexity of Interior-Point Methods

## General Complexity

Complexity for LP, QP, QCQP, SOCP, and SDP is  $O(n^3L)$ , where

- $n$ : number of variables.
- $L$ : number of accuracy digits of the solution.

## Specific Complexities

- **LP**:  $O((m+n)^{3/2}n^2L)$ .
- **QCQP**:  $O(\sqrt{m}(m+n)n^2L)$ .
- **SOCP**:  $O(\sqrt{m+1}n(n^2+m+(m+1)k^2)L)$  with  $k$  the cone dimension.
- **SDP**:  $O(\sqrt{1+mk}n(n^2+nmk^2+mk^3)L)$ , with  $k \times k$  the matrix dimension.

## Example Analysis

- For SOCP with  $m = O(n)$  and  $k = O(n)$ , complexity is  $O(n^{4.5}L)$ .
- For SDP with  $k = O(n)$ , complexity is  $O(n^4)$ .
- If  $m = O(n)$  for SDP, complexity becomes  $O(n^6L)$ .
- Complexity for solving SOCP is higher than LP, QP, and QCQP; even higher for SDP.



## Solvers and Standard Form:

- Problems must be expressed in a standard form for solvers.
- Transformation to standard form is time-consuming and error-prone.

## General Norm Approximation Problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{Ax} - \mathbf{b}\|$$

Solution depends on the choice of the norm.

## Norm Approximation with Euclidean or $\ell_2$ -norm:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{Ax} - \mathbf{b}\|_2$$

Least squares (LS) problem with analytic solution:  $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ .

# Interface with Solvers

## Norm Approximation with Chebyshev or $\ell_\infty$ -norm:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{Ax} - \mathbf{b}\|_\infty$$

- Rewritten as LP:

$$\begin{aligned} &\underset{\mathbf{x}, t}{\text{minimize}} && t \\ &\text{subject to} && -t\mathbf{1} \leq \mathbf{Ax} - \mathbf{b} \leq t\mathbf{1} \end{aligned}$$

- Equivalent form:

$$\begin{aligned} &\underset{\mathbf{x}, t}{\text{minimize}} && \begin{bmatrix} \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ t \end{bmatrix} \\ &\text{subject to} && \begin{bmatrix} \mathbf{A} & -\mathbf{1} \\ -\mathbf{A} & -\mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ t \end{bmatrix} \leq \begin{bmatrix} \mathbf{b} \\ -\mathbf{b} \end{bmatrix} \end{aligned}$$

- Matlab code:

```
xt = linprog( [zeros(n,1); 1],  
             [A,-ones(m,1); -A,-ones(m,1)],  
             [b; -b] )  
x = xt(1:n)
```

## Norm Approximation Problem with Manhattan or $\ell_1$ -Norm:

- Rewritten as LP:
$$\begin{aligned} & \underset{x}{\text{minimize}} && \|Ax - b\|_1 \\ & \underset{x,t}{\text{minimize}} && \mathbf{1}^T t \\ & \text{subject to} && -t \leq Ax - b \leq t \end{aligned}$$
- Equivalent form:

$$\begin{aligned} & \underset{x,t}{\text{minimize}} && \begin{bmatrix} \mathbf{0}^T & \mathbf{1}^T \end{bmatrix} \begin{bmatrix} x \\ t \end{bmatrix} \\ & \text{subject to} && \begin{bmatrix} A & -I \\ -A & -I \end{bmatrix} \begin{bmatrix} x \\ t \end{bmatrix} \leq \begin{bmatrix} b \\ -b \end{bmatrix} \end{aligned}$$

- Matlab code:

```
xt = linprog( [zeros(n,1); ones(n,1)],  
             [A,-eye(m,1); -A,-eye(m,1)],  
             [b; -b] )  
  
x = xt(1:n)
```

## Euclidean Norm Approximation Problem with Linear Constraints:

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & \|\mathbf{Ax} - \mathbf{b}\|_2 \\ \text{subject to} & \mathbf{Cx} = \mathbf{d} \\ & \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}.\end{array}$$

- Equivalent form:

$$\begin{array}{ll}\underset{\mathbf{x}, \mathbf{y}, t, \mathbf{s}_l, \mathbf{s}_u}{\text{minimize}} & t \\ \text{subject to} & \mathbf{Ax} - \mathbf{b} = \mathbf{y} \\ & \mathbf{Cx} = \mathbf{d} \\ & \mathbf{x} - \mathbf{s}_l = \mathbf{l} \\ & \mathbf{x} + \mathbf{s}_u = \mathbf{u} \\ & \mathbf{s}_l, \mathbf{s}_u \geq \mathbf{0} \\ & \|\mathbf{y}\|_2 \leq t\end{array}$$

# Interface with Solvers

- Recall the previous equivalent form:

$$\begin{array}{ll}\text{minimize} & t \\ \text{subject to} & \mathbf{Ax} - \mathbf{b} = \mathbf{y} \\ & \mathbf{Cx} = \mathbf{d} \\ & \mathbf{x} - \mathbf{s}_l = \mathbf{l} \\ & \mathbf{x} + \mathbf{s}_u = \mathbf{u} \\ & \mathbf{s}_l, \mathbf{s}_u \geq \mathbf{0} \\ & \|\mathbf{y}\|_2 \leq t\end{array}$$

- Equivalent form in matrix notation:

$$\begin{array}{ll}\text{minimize} & \begin{bmatrix} \mathbf{0}^T & \mathbf{0}^T & \mathbf{0}^T & \mathbf{0}^T & 1 \end{bmatrix} \bar{\mathbf{x}} \\ \text{subject to} & \begin{bmatrix} \mathbf{A} & & & -\mathbf{I} \\ \mathbf{C} & & & \\ \mathbf{I} & -\mathbf{I} & & \\ \mathbf{I} & & \mathbf{I} & \end{bmatrix} \bar{\mathbf{x}} \leq \begin{bmatrix} \mathbf{b} \\ \mathbf{d} \\ \mathbf{l} \\ \mathbf{u} \end{bmatrix} \\ & \bar{\mathbf{x}} \in \mathbf{R}^n \times \mathbf{R}_+^n \times \mathbf{R}_+^n \times \mathbf{Q}^m\end{array}$$

- Matlab code:

```
AA = [ A, zeros(m,n), zeros(m,n), -eye(m), 0;  
       C, zeros(p,n), zeros(p,n), zeros(p,n), 0;  
       eye(n), -eye(n), zeros(n,n), zeros(n,n), 0;  
       eye(n), zeros(n,n), eye(n), zeros(n,n), 0 ]  
bb = [ b; d; l; u ]  
cc = [ zeros(3*n + m, 1); 1 ]  
K.f = n; K.l = 2*n; K.q = m + 1  
xsyz = sedumi( AA, bb, cc, K )  
x = xsyz(1:n)
```

## Modeling Framework:

- Simplifies the use of solvers by handling solver argument formatting.
- Acts as an interface between the user and the solver.
- Can interface with various solvers, allowing user choice based on problem type.
- Useful for rapid prototyping and avoiding transcription errors.
- Direct solver calls may be preferred for high-speed requirements.

## Successful Examples:

- YALMIP: For Matlab (Löfberg 2004).
- CVX: Initially released in 2005 for Matlab. Now available in Python, R, and Julia. (Grant and Boyd 2008, 2014; Fu, Narasimhan, and Boyd 2020).

## CVX (Convex Disciplined Programming):

- Tool for rapid prototyping of models and algorithms with convex optimization.
- Interfaces with solvers like SeDuMi, SDPT3, Gurobi, and MOSEK.
- Recognizes elementary convex and concave functions and composition rules.
- Determines problem convexity.
- Simple and convenient for prototyping.

## Example: Constrained Euclidean Norm Approximation in CVX:

- Problem statement:

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & \|\mathbf{Ax} - \mathbf{b}\|_2 \\ \text{subject to} & \mathbf{Cx} = \mathbf{d} \\ & \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}\end{array}$$

- Matlab code:

```
cvx_begin
    variable x(n)
    minimize(norm(A * x - b, 2))
    subject to
        C * x == d
        l <= x
        x <= u
cvx_end
```



## Example: Constrained Euclidean Norm Approximation in CVX: (cont'd)

- R code:

```
x <- Variable(n)
prob <- Problem(Minimize(cvxr_norm(A %*% x - b, 2)),
               list(C %*% x == d,
                    l <= x,
                    x <= u))

solve(prob)
```

- Python code:

```
x = cvxpy.Variable(n)
prob = cvxpy.Problem(cvxpy.Minimize(cvxpy.norm(A @ x - b, 2)),
                    [C @ x == d,
                     l <= x,
                     x <= u])

prob.solve()
```

# Outline

- 1 Solvers
- 2 Gradient Methods**
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary

## Unconstrained Optimization Problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x})$$

where  $f$  is the objective function, assumed to be continuously differentiable.

## Iterative Methods:

- Produce a sequence of iterates  $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$
- Sequence may or may not converge to an optimal solution  $\mathbf{x}^*$ .

**Ideal Case with Convex  $f$ :** As iterations proceed ( $k \rightarrow \infty$ ),

- objective function converges to the optimal value:

$$f(\mathbf{x}^k) \rightarrow p^*$$

- gradient tends to zero:

$$\nabla f(\mathbf{x}^k) \rightarrow \mathbf{0}$$

**References:** (Bertsekas 1999), (S. P. Boyd and Vandenberghe 2004), (Nocedal and Wright 2006), (Beck 2017).

## Descent Methods (Gradient Methods):

- Satisfy the property:  $f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k)$ .
- Iterates are obtained as:

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k,$$

- $\mathbf{d}^k$ : *search direction*.
- $\alpha^k$ : *stepsize* at iteration  $k$ .

## Descent Property:

- For a sufficiently small step,  $\mathbf{d}$  must satisfy:

$$\nabla f(\mathbf{x})^\top \mathbf{d} < 0$$

- $\alpha$  must be properly chosen (if too large, the descent property may be violated even with a descent direction).

**Line Search:** Procedure to choose the stepsize  $\alpha$ .

- **Exact Line Search:** Solves the univariate optimization problem

$$\alpha = \arg \min_{\alpha > 0} f(\mathbf{x} + \alpha \mathbf{d}).$$

- **Backtracking Line Search (Armijo Rule):** Starting at  $\alpha = 1$ , repeat  $\alpha \leftarrow \beta \alpha$  until

$$f(\mathbf{x} + \alpha \mathbf{d}) < f(\mathbf{x}) + \sigma \alpha \nabla f(\mathbf{x})^T \mathbf{d},$$

where  $\sigma \in (0, 1/2)$  and  $\beta \in (0, 1)$  are given parameters.

# Gradient Descent Method

## Gradient Descent Method (Steepest Descent Method)

A descent method where the search direction is the opposite of the gradient:

$$\mathbf{d} = -\nabla f(\mathbf{x}),$$

which is a descent direction since  $\nabla f(\mathbf{x})^\top \mathbf{d} < 0$ .

## Gradient Descent Update:

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

**Stopping Criterion:** Common heuristic:  $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$ .

## Convergence:

- Often slow, making it rarely used in practice.
- Useful in high-dimensional problems or when distributed implementation is required.

## Gradient descent method

### Initialization:

- Choose initial point  $\mathbf{x}^0$ .
- Set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Compute the negative gradient as descent direction:  $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$ .
- 2 Line search: Choose a stepsize  $\alpha^k > 0$  via exact or backtracking line search.
- 3 Obtain next iterate:

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

- 4  $k \leftarrow k + 1$

**Until:** convergence

## Newton's Method

- A descent method using both the gradient and the Hessian of  $f$ .
- **Search direction:**

$$\mathbf{d} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}),$$

- Assumes  $f$  is convex, twice continuously differentiable, and the Hessian matrix is positive definite for all  $\mathbf{x}$ .

**Second-Order Approximation:**  $\mathbf{x} + \mathbf{d}$  minimizes the second-order approximation of  $f(\mathbf{x})$  around  $\mathbf{x}$ :

$$\hat{f}(\mathbf{x} + \mathbf{v}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v}.$$

## Newton's Method Update:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha^k \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k).$$



## Newton Decrement

- Measures the proximity of  $\mathbf{x}$  to an optimal point:

$$\lambda(\mathbf{x}) = (\nabla f(\mathbf{x})^T \nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}))^{1/2}$$

- Estimates  $f(\mathbf{x}) - p^*$ :

$$f(\mathbf{x}) - \inf_{\mathbf{y}} \hat{f}(\mathbf{y}) = \frac{1}{2} \lambda(\mathbf{x})^2$$

- Computational cost of the Newton decrement is negligible since  $\lambda(\mathbf{x})^2 = -\nabla f(\mathbf{x})^T \mathbf{d}$ .

## Advantages and Limitations

- Fast convergence.
- Central to most modern solvers.
- Impractical for very large dimensional problems due to computation and storage of the Hessian.
- For large problems, quasi-Newton methods are used (Nocedal and Wright 2006).

## Newton's method

### Initialization:

- Choose initial point  $\mathbf{x}^0$  and tolerance  $\epsilon > 0$ . Set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Compute Newton direction and decrement:

$$\mathbf{d}^k = -\nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k) \quad \text{and} \quad \lambda(\mathbf{x}^k)^2 = -\nabla f(\mathbf{x}^k)^\top \mathbf{d}^k.$$

- 2 Line search: Choose a stepsize  $\alpha^k > 0$  via exact or backtracking line search.
- 3 Obtain next iterate:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha^k \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k).$$

- 4  $k \leftarrow k + 1$

**Until:** convergence (i.e.,  $\lambda(\mathbf{x}^k)^2/2 \leq \epsilon$ )

## Convergence of Descent Methods

- Ideally, the sequence  $\{\mathbf{x}^k\}$  should converge to a global minimum.
- For non-convex  $f$ , convergence to a global minimum is unlikely due to local minima.
- Descent methods typically converge to a stationary point.
- For convex  $f$ , a stationary point is a global minimum.

## Theoretical Convergence

- Descent methods have nice theoretical convergence properties (Bertsekas 1999).
- **Theorem: Convergence of descent methods**
  - Suppose  $\{\mathbf{x}^k\}$  is a sequence generated by a descent method (e.g., gradient descent or Newton's method).
  - Stepsize  $\alpha^k$  chosen by exact line search or backtracking line search.
  - Every limit point of  $\{\mathbf{x}^k\}$  is a stationary point of the problem.

## Simpler Stepsize Rules with Theoretical Convergence (Bertsekas 1999)

- Constant stepsize:  $\alpha^k = \alpha$  for sufficiently small  $\alpha$ .
- Diminishing stepsize rule:  $\alpha^k \rightarrow 0$  with  $\sum_{k=0}^{\infty} \alpha^k = \infty$ .

## Newton's Method Convergence Phases

- Damped Newton phase: Slow convergence.
- Quadratically convergent phase: Extremely fast convergence.

## Practical Considerations

- Gradient descent method converges slowly.
- Newton's method converges much faster but requires computing the Hessian.
- Newton's method is preferred if problem dimensionality is manageable.
- For extremely large dimensional problems (e.g., deep learning), computing and storing the Hessian is not feasible.

# Projected Gradient Methods

## Constrained Optimization Problem

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{X},\end{array}$$

where  $f$  is the objective function (continuously differentiable) and  $\mathcal{X}$  is a convex set.

## Descent Method

- Iterative update:

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$$

where  $\mathbf{d}^k$  is a descent direction.

- Potential issue:  $\mathbf{x}^{k+1}$  may be infeasible.

## Projected Gradient Methods (Gradient Projection Methods)

Address infeasibility by projecting onto the feasible set after taking the step (Bertsekas 1999; Beck 2017):

$$\mathbf{x}^{k+1} = \left[ \mathbf{x}^k + \alpha^k \mathbf{d}^k \right]_{\mathcal{X}}$$

where  $[\mathbf{x}]_{\mathcal{X}}$  denotes projection onto the set  $\mathcal{X}$ :  $\min_{\mathbf{y}} \|\mathbf{y} - \mathbf{x}\|$  subject to  $\mathbf{y} \in \mathcal{X}$ .

# Projected Gradient Methods

## Generalized Gradient Projection Method

- Iterative update:

$$\begin{aligned}\bar{\mathbf{x}}^k &= \left[ \mathbf{x}^k + s^k \mathbf{d}^k \right]_{\mathcal{X}} \\ \mathbf{x}^{k+1} &= \mathbf{x}^k + \alpha^k \left( \bar{\mathbf{x}}^k - \mathbf{x}^k \right),\end{aligned}$$

- $\mathbf{d}^k = \bar{\mathbf{x}}^k - \mathbf{x}^k$  is a feasible direction.
- $\alpha^k$  is the stepsize.
- $s^k$  is a positive scalar (Bertsekas 1999).
- Special case:  $\alpha^k = 1$ :

$$\mathbf{x}^{k+1} = \left[ \mathbf{x}^k + s^k \mathbf{d}^k \right]_{\mathcal{X}}$$

- $s^k$  can be viewed as a stepsize.
- If  $\mathbf{x}^k + s^k \mathbf{d}^k$  is already feasible, the method reduces to the regular gradient method.

## Practical Consideration

Gradient projection method is practical only if the projection is easy to compute.

## Projected Gradient Descent Method

- Uses the negative gradient as the search direction.
- Iterative update:

$$\begin{aligned}\bar{\mathbf{x}}^k &= \left[ \mathbf{x}^k - s^k \nabla f(\mathbf{x}^k) \right]_{\mathcal{X}} \\ \mathbf{x}^{k+1} &= \mathbf{x}^k + \alpha^k (\bar{\mathbf{x}}^k - \mathbf{x}^k),\end{aligned}$$

- $\bar{\mathbf{x}}^k$ : Projection of  $\mathbf{x}^k - s^k \nabla f(\mathbf{x}^k)$  onto the set  $\mathcal{X}$ .
- $\alpha^k$ : Stepsize.
- $s^k$ : Positive scalar stepsize for the gradient step.

## Constrained Newton's Method

- Assumptions:
  - $f$  is twice continuously differentiable.
  - The Hessian matrix is positive definite for all  $\mathbf{x} \in \mathcal{X}$ .
- Iterative update:

$$\bar{\mathbf{x}}^k = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \nabla f(\mathbf{x}^k)^\top (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2s^k} (\mathbf{x} - \mathbf{x}^k)^\top \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k) \right\}$$
$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k (\bar{\mathbf{x}}^k - \mathbf{x}^k).$$

- $\bar{\mathbf{x}}^k$ : Solution to the quadratic subproblem.
  - $\alpha^k$ : Stepsize.
  - $s^k$ : Positive scalar.
- Observations:
  - If  $s^k = 1$ , the quadratic cost is the second-order Taylor series expansion of  $f$  around  $\mathbf{x}^k$ .
  - The main difficulty is solving the quadratic subproblem to find  $\bar{\mathbf{x}}^k$ .
    - This may not be simple even when the constraint set  $\mathcal{X}$  has a simple structure.
    - The method typically makes practical sense only for problems of small dimension.



## Convergence of Gradient Projection Methods

- Detailed in (Bertsekas 1999).
- **Theorem: Convergence of gradient projection methods.**
  - Suppose  $\{\mathbf{x}^k\}$  is a sequence generated by a gradient projection method (e.g., projected gradient descent method or constrained Newton's method).
  - Stepsize  $\alpha^k$  chosen by exact line search or backtracking line search.
  - Every limit point of  $\{\mathbf{x}^k\}$  is a stationary point of the problem.

## Simpler Stepsize Rules with Theoretical Convergence

- Constant stepsize:  $\alpha^k = 1$  and  $s^k = s$  for sufficiently small  $s$  (Bertsekas 1999).

# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)**
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary

# Interior-Point Methods (IPM)

## Traditional Optimization Algorithms

- Based on gradient projection methods.
- May suffer from:
  - Slow convergence.
  - Sensitivity to algorithm initialization.
  - Sensitivity to stepsize selection.

## Interior-Point Methods (IPM)

- Modern approach for convex problems.
- Enjoy excellent convergence properties (polynomial convergence).
- Do not suffer from the usual problems of traditional methods.

## Convex Optimization Problem

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & \mathbf{Ax} = \mathbf{b}\end{array}$$

**References:** (Nesterov and Nemirovskii 1994; Bertsekas 1999; Nemirovski 2001; S. P. Boyd and Vandenberghe 2004; Nocedal and Wright 2006; Nesterov 2018).

# Eliminating Equality Constraints

## Dealing with Equality Constraints

- Can be handled via Lagrange duality (S. P. Boyd and Vandenberghe 2004).
- Alternatively, can be eliminated in a pre-processing stage.

## Representation of Solutions to $\mathbf{Ax} = \mathbf{b}$ :

$$\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{Ax} = \mathbf{b}\} = \{\mathbf{Fz} + \mathbf{x}_0 \mid \mathbf{z} \in \mathbb{R}^{n-p}\},$$

- $\mathbf{x}_0$  is any particular solution to  $\mathbf{Ax} = \mathbf{b}$ .
- $\mathbf{F} \in \mathbb{R}^{n \times (n-p)}$  spans the nullspace of  $\mathbf{A}$ , i.e.,  $\mathbf{AF} = \mathbf{0}$ .

## Reduced or Eliminated Problem

- Equivalent to the original problem:

$$\begin{aligned} & \underset{\mathbf{z}}{\text{minimize}} && \tilde{f}_0(\mathbf{z}) \triangleq f_0(\mathbf{Fz} + \mathbf{x}_0) \\ & \text{subject to} && \tilde{f}_i(\mathbf{z}) \triangleq f_i(\mathbf{Fz} + \mathbf{x}_0) \leq 0, \quad i = 1, \dots, m, \end{aligned}$$

- Gradients and Hessians:

$$\begin{aligned} \nabla \tilde{f}_i(\mathbf{z}) &= \mathbf{F}^T \nabla f_i(\mathbf{x}) \\ \nabla^2 \tilde{f}_i(\mathbf{z}) &= \mathbf{F}^T \nabla^2 f_i(\mathbf{x}) \mathbf{F}. \end{aligned}$$

## Reformulation via Indicator Function:

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & f_0(\mathbf{x}) + \sum_{i=1}^m l_-(f_i(\mathbf{x})) \\ \text{subject to} & \mathbf{Ax} = \mathbf{b}, \end{array}$$

using the indicator function:

$$l_-(u) = \begin{cases} 0 & \text{if } u \leq 0 \\ \infty & \text{otherwise.} \end{cases}$$

## Characteristics of the Reformulated Problem

- The inequality constraints are eliminated.
- The indicator function is included in the objective.
- The main drawbacks are that the indicator function is noncontinuous and nondifferentiable, making the approach not practical for optimization.

## Logarithmic Barrier

A popular smooth approximation of the indicator function is the **logarithmic barrier**:

$$l_-(u) \approx -\frac{1}{t} \log(-u),$$

- The parameter  $t > 0$  controls the approximation.
- The approximation improves as  $t \rightarrow \infty$ .

## Approximate Problem Using the Logarithmic Barrier

The reformulated problem is:

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & f_0(\mathbf{x}) - \frac{1}{t} \sum_{i=1}^m \log(-f_i(\mathbf{x})) \\ \text{subject to} & \mathbf{Ax} = \mathbf{b}. \end{array}$$

# Logarithmic Barrier Method

## Logarithmic Barrier Function

- The overall barrier function (excluding the  $1/t$  factor) is:

$$\phi(\mathbf{x}) = - \sum_{i=1}^m \log(-f_i(\mathbf{x})),$$

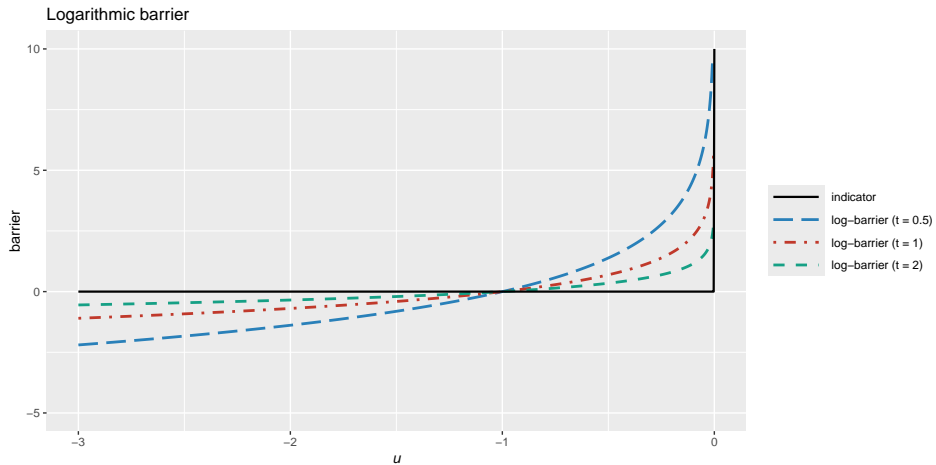
which is convex (from composition rules).

- The gradient and Hessian are:

$$\begin{aligned}\nabla \phi(\mathbf{x}) &= \sum_{i=1}^m \frac{1}{-f_i(\mathbf{x})} \nabla f_i(\mathbf{x}) \\ \nabla^2 \phi(\mathbf{x}) &= \sum_{i=1}^m \frac{1}{f_i(\mathbf{x})^2} \nabla f_i(\mathbf{x}) \nabla f_i(\mathbf{x})^\top + \sum_{i=1}^m \frac{1}{-f_i(\mathbf{x})} \nabla^2 f_i(\mathbf{x}).\end{aligned}$$

# Logarithmic Barrier Method

Logarithmic barrier for several values of the parameter  $t$ :





# Central Path

## Central Path

Defined as the curve  $\{\mathbf{x}^*(t) \mid t > 0\}$ , where  $\mathbf{x}^*(t)$  is the solution to

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & tf_0(\mathbf{x}) + \phi(\mathbf{x}) \\ \text{subject to} & \mathbf{Ax} = \mathbf{b}, \end{array}$$

which can be solved via Newton's method.

## Solution to the Central Path Problem

- Ignoring equality constraints for simplicity:

$$t\nabla f_0(\mathbf{x}) + \sum_{i=1}^m \frac{1}{-f_i(\mathbf{x})} \nabla f_i(\mathbf{x}) = \mathbf{0}.$$

- Define  $\lambda_i^*(t) = 1/(-tf_i(\mathbf{x}^*(t)))$ .
- $\mathbf{x}^*(t)$  minimizes the Lagrangian:

$$L(\mathbf{x}; \boldsymbol{\lambda}^*(t)) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i^*(t) f_i(\mathbf{x}).$$

## Convergence to Optimal Value

As  $t \rightarrow \infty$ ,  $f_0(\mathbf{x}^*(t)) \rightarrow p^*$ . From Lagrange duality theory:

$$\begin{aligned} p^* &\geq g(\boldsymbol{\lambda}^*(t)) \\ &= L(\mathbf{x}^*(t); \boldsymbol{\lambda}^*(t)) \\ &= f_0(\mathbf{x}^*(t)) - m/t. \end{aligned}$$

## Connection with KKT Conditions

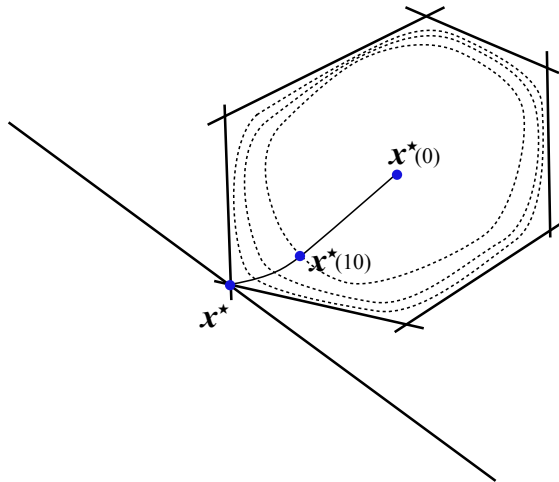
- The pair  $\mathbf{x}^*(t)$  and  $\boldsymbol{\lambda}^*(t)$  satisfies:

$$\begin{array}{ll} f_i(\mathbf{x}) \leq 0, & i = 1, \dots, m \quad (\text{primal feasibility}) \\ \lambda_i \geq 0, & i = 1, \dots, m \quad (\text{dual feasibility}) \\ \lambda_i f_i(\mathbf{x}) = -\frac{1}{t}, & i = 1, \dots, m \quad (\text{approximate complementary slackness}) \\ \nabla f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i \nabla f_i(\mathbf{x}) = \mathbf{0}. & (\text{zero Lagrangian gradient}) \end{array}$$

- The difference with the original KKT conditions is:
  - Complementary slackness is approximately satisfied.
  - The approximation improves as  $t \rightarrow \infty$ .

# Central Path

Example of central path of an LP:



# Logarithmic Barrier Method

## Smooth Approximation with Logarithmic Barrier

- The log-barrier problem is a smooth approximation of the original problem.
- The approximation improves as  $t \rightarrow \infty$ .

## Challenges with Choosing $t$

- Large  $t$ :
  - Leads to slow convergence.
  - Gradients and Hessians vary greatly near the boundary of the feasible set.
  - Newton's method fails to reach quadratic convergence.
- Small  $t$ :
  - Facilitates better convergence.
  - The approximation is not close to the original problem.

## Adaptive $t$ Approach

- Change  $t$  over iterations to balance fast convergence and accurate approximation.
- At each outer iteration, update  $t$  and compute  $\mathbf{x}^*(t)$  using Newton's method.
- Interior-point methods (IPM):
  - Achieve this trade-off.
  - For each  $t > 0$ ,  $\mathbf{x}^*(t)$  is strictly feasible and lies in the interior of the feasible set.

## Barrier Method

- A type of primal-based IPM.
- Update rule for  $t$ :
  - $t^{k+1} \leftarrow \mu t^k$ , where  $\mu > 1$ .
  - Typically,  $t^0 = 1$ .
- Choice of  $\mu$ :
  - Large  $\mu$  means fewer outer iterations but more inner (Newton) iterations.
  - Typical values:  $\mu = 10 \sim 20$ .
- Termination criterion:
  - $m/t < \epsilon$ , guaranteeing  $f_0(\mathbf{x}) - p^* \leq \epsilon$ .
- Refer to (S. P. Boyd and Vandenberghe 2004) for practical details.

## Logarithmic barrier method for constrained optimization

### Initialization:

- Choose initial point  $\mathbf{x}^0 \in \mathcal{X}$  strictly feasible,  $t^0 > 0$ ,  $\mu > 1$ , and tolerance  $\epsilon > 0$ .
- Set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Centering step: compute next iterate  $\mathbf{x}^{k+1}$  by solving the central path problem with  $t = t^k$  and initial point  $\mathbf{x}^k$ .
- 2 Increase  $t$ :  $t^{k+1} \leftarrow \mu t^k$ .
- 3  $k \leftarrow k + 1$

**Until:** convergence (i.e.,  $m/t < \epsilon$ )

# Logarithmic Barrier Method

## Example: Barrier Method for LP

Consider the LP:

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & \mathbf{c}^T \mathbf{x} \\ \text{subject to} & \mathbf{Ax} \leq \mathbf{b}.\end{array}$$

## Convergence Analysis

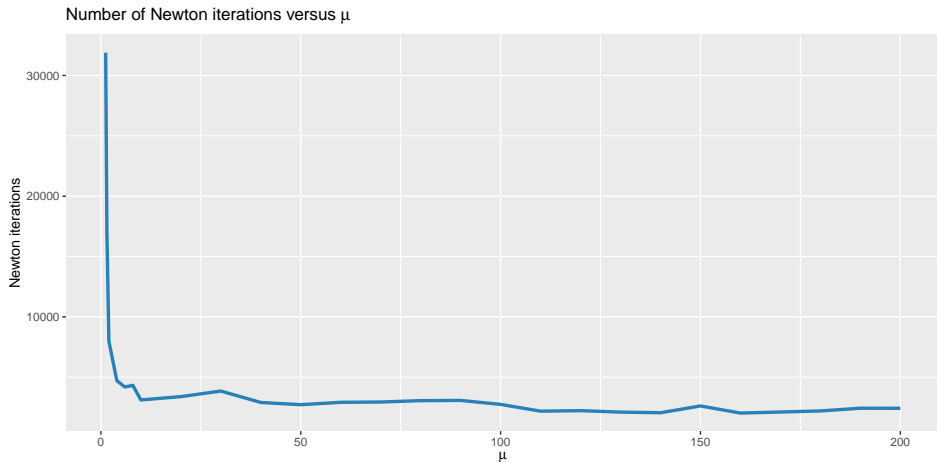
- Use the barrier method with different  $\mu$  values.
- The case is  $m = 100$  inequalities and  $n = 50$  variables.
- The duality gap is  $\epsilon = 10^{-6}$ .
- The centering problem is solved via Newton's method.

## Observation

The total number of Newton iterations is not very sensitive to  $\mu$  as long as  $\mu \geq 10$ .

# Logarithmic Barrier Method

Convergence of barrier method for an LP for different values of  $\mu$ :





# Convergence

## Termination Criterion

- Number of outer iterations (centering steps) required:

$$\frac{m}{\mu^k t^0} \leq \epsilon$$

- Solving for  $k$ :

$$\left\lceil \frac{\log(m / (\epsilon t^0))}{\log(\mu)} \right\rceil$$

where  $\lceil \cdot \rceil$  is the ceiling operator.

## Convergence of Centering Steps

- Characterized via the convergence for Newton's method.
- Specific updates for  $\mu$  and good initialization points for each centering step are not considered in this simple analysis.

## References

For detailed convergence analysis: (Nesterov and Nemirovskii 1994; Nemirovski 2001; S. P. Boyd and Vandenberghe 2004; Nocedal and Wright 2006; Nesterov 2018).

## Barrier Method and Strictly Feasible Initial Point

- The barrier method requires a strictly feasible initial point  $\mathbf{x}^0$  (such that  $f_i(\mathbf{x}^0) < 0$ ).
- If such a point is not known, a preliminary stage (*Phase I*) is used to find it.
- The barrier method itself is then called *Phase II*.

## Phase I Methods

- Aim to find a feasible point for the original problem by solving the feasibility problem:

$$\begin{array}{ll}\text{find} & \mathbf{x} \\ \text{subject to} & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & \mathbf{Ax} = \mathbf{b}.\end{array}$$

- The barrier method cannot be used directly for the feasibility problem as it requires a feasible starting point.

# Feasibility and Phase I Methods

## Formulating Phase I Methods

A simple example involves solving the convex optimization problem:

$$\begin{array}{ll}\underset{\mathbf{x}, s}{\text{minimize}} & s \\ \text{subject to} & f_i(\mathbf{x}) \leq s, \quad i = 1, \dots, m \\ & \mathbf{Ax} = \mathbf{b}.\end{array}$$

- To construct a strictly feasible point, choose any  $\mathbf{x}$  that satisfies the equality constraints.
- Then, choose  $s$  such that  $s > f_i(\mathbf{x})$ , e.g.,  $s = 1.1 \times \max_i \{f_i(\mathbf{x})\}$ .
- This provides an initial strictly feasible point for the Phase I problem.

## Solving the Phase I Problem

- After obtaining  $(\mathbf{x}^*, s^*)$ , the value of  $s^*$  is checked.
- If  $s^* < 0$ , then  $\mathbf{x}^*$  is a strictly feasible point and can be used in the barrier method to solve the original problem.
- If  $s^* > 0$ , then no feasible point exists, and there is no need to attempt solving the original problem as it is infeasible.

# Primal-Dual Interior-Point Methods

## Primal Barrier Method

- Requires a strictly feasible initial point.
- Involves distinct inner and outer iterations.

## Primal-Dual IPMs

- More efficient, especially for high accuracy.
- Exhibit superlinear asymptotic convergence.
- Key features:
  - Update both primal and dual variables at each iteration.
  - No distinction between inner and outer iterations.
  - Can start at infeasible points, eliminating the need for Phase I methods.

## Advantages

- Efficiency: Better for high accuracy.
- Convergence: Superlinear asymptotic convergence.
- Initialization: Can start from infeasible points, simplifying the process.

**Summary:** Primal-dual IPMs offer significant advantages over the primal barrier method in terms of efficiency, convergence, and ease of initialization.

# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods**
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary

# Fractional Programming (FP) Methods

## Concave-Convex Fractional Program (FP):

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{maximize}} & \frac{f(\mathbf{x})}{g(\mathbf{x})} \\ \text{subject to} & \mathbf{x} \in \mathcal{X},\end{array}$$

- $f(\mathbf{x})$  is a concave function
- $g(\mathbf{x}) > 0$  is a convex function
- $\mathcal{X}$  is a convex feasible set

## Nature of Fractional Programs

- Nonconvex problems, generally difficult to solve.
- Concave-convex FP is a quasiconvex optimization problem, making it more tractable.

## Methods to Solve Concave-Convex Fractional Programs:

- Iterative bisection method
- Dinkelbach method
- Schaible transform

## Problem Reformulation

- Solve a sequence of convex feasibility problems:

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{find}} & \mathbf{x} \\ \text{subject to} & tg(\mathbf{x}) \leq f(\mathbf{x}) \\ & \mathbf{x} \in \mathcal{X} \end{array}$$

- $t > 0$  is a fixed parameter, not an optimization variable.

**Goal:** Find the optimal value of  $t$  for the original problem.

## Procedure:

- If the feasibility problem is infeasible,  $t$  is too large and must be decreased.
- If feasible,  $t$  is too small and can be increased.

## Key points:

- Starts with an interval  $[l, u]$  known to contain the optimal value  $p^*$  and sequentially halves the interval.
- The length of the interval after  $k$  iterations is  $2^{-k}(u - l)$ .
- Number of iterations required to achieve a tolerance of  $\epsilon$  is  $\lceil \log_2((u - l)/\epsilon) \rceil$ .

## Bisection method (aka “sandwich technique”) for concave-convex FP

### Initialization:

- Initialize  $l$  and  $u$  such that  $p^* \in [l, u]$ .

### Repeat while $(u - l) > \epsilon$ :

- Compute midpoint of interval:  $t = (l + u)/2$ .
- Solve the convex feasibility problem for  $t$ .
- If feasible, set  $u = t$ ; otherwise set  $l = t$ .



## Dinkelbach Transform

- Objective: reformulate the original concave-convex FP into a sequence of simpler convex problems.
- Reformulated problem:

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{maximize}} & f(\mathbf{x}) - y^k g(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{X}\end{array}$$

- Parameter update:  $y^k = \frac{f(\mathbf{x}^k)}{g(\mathbf{x}^k)}$  with  $k$  as the iteration index.

## Convergence

- The Dinkelbach method converges to global optimum of the original concave-convex FP.
- Key properties:
  - Increasing sequence  $\{y^k\}$ .
  - Function  $F(y) = \arg \max_{\mathbf{x}} \{f(\mathbf{x}) - yg(\mathbf{x})\}$ .

## Key points:

- Transforms a nonconvex problem into a sequence of convex problems.
- Ensures global optimality through iterative updates.

## Dinkelback method for concave-convex FP

### Initialization:

- Choose initial point  $\mathbf{x}^0$ .
- Set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Set  $y^k = f(\mathbf{x}^k)/g(\mathbf{x}^k)$ .
- 2 Solve the reformulated convex problem and keep current solution as  $\mathbf{x}^{k+1}$ .
- 3  $k \leftarrow k + 1$

**Until:** convergence

# FP: Charnes-Cooper Transform

## Linear Fractional Program (LFP)

$$\begin{array}{ll}\text{minimize}_{\mathbf{x}} & \frac{\mathbf{c}^T \mathbf{x} + d}{\mathbf{e}^T \mathbf{x} + f} \\ \text{subject to} & \mathbf{G}\mathbf{x} \leq \mathbf{h} \\ & \mathbf{A}\mathbf{x} = \mathbf{b}\end{array}$$

$$\text{with dom } f_0 = \left\{ \mathbf{x} \mid \mathbf{e}^T \mathbf{x} + f > 0 \right\}.$$

**Charnes-Cooper Transform:** Transforms the original LFP into a linear program (LP):

$$\begin{array}{ll}\text{minimize}_{\mathbf{y}, t} & \mathbf{c}^T \mathbf{y} + dt \\ \text{subject to} & \mathbf{G}\mathbf{y} \leq \mathbf{h}t \\ & \mathbf{A}\mathbf{y} = \mathbf{b}t \\ & \mathbf{e}^T \mathbf{y} + ft = 1 \\ & t \geq 0\end{array}$$

$$\text{where } \mathbf{y} = \frac{\mathbf{x}}{\mathbf{e}^T \mathbf{x} + f} \text{ and } t = \frac{1}{\mathbf{e}^T \mathbf{x} + f}.$$

## Proof:

- Any feasible point  $\mathbf{x}$  in the original LFP leads to a feasible point  $(\mathbf{y}, t)$  in the LP with the same objective value.
- Conversely, any feasible point  $(\mathbf{y}, t)$  in the LP leads to a feasible point  $\mathbf{x}$  in the original LFP via  $\mathbf{x} = \mathbf{y}/t$ , also with the same objective value:

$$\frac{\mathbf{c}^T \mathbf{y} + dt}{1} = \frac{\mathbf{c}^T \mathbf{y} + dt}{\mathbf{e}^T \mathbf{y} + ft} = \frac{\mathbf{c}^T \mathbf{y}/t + d}{\mathbf{e}^T \mathbf{y}/t + f} = \frac{\mathbf{c}^T \mathbf{x} + d}{\mathbf{e}^T \mathbf{x} + f}.$$

## Concave-Convex Fractional Program (FP)

$$\begin{array}{ll}\text{maximize}_{\mathbf{x}} & \frac{f(\mathbf{x})}{g(\mathbf{x})} \\ \text{subject to} & \mathbf{x} \in \mathcal{X}\end{array}$$

**Schaible Transform:** Rewrites the original concave-convex FP into a convex problem:

$$\begin{array}{ll}\text{maximize}_{\mathbf{y}, t} & tf\left(\frac{\mathbf{y}}{t}\right) \\ \text{subject to} & tg\left(\frac{\mathbf{y}}{t}\right) \leq 1 \\ & t \geq 0 \\ & \mathbf{y}/t \in \mathcal{X}\end{array}$$

where  $\mathbf{y} = \frac{\mathbf{x}}{g(\mathbf{x})}$  and  $t = \frac{1}{g(\mathbf{x})}$ .

## Proof:

- Any feasible point  $\mathbf{x}$  in the original FP leads to a feasible point  $(\mathbf{y}, t)$  in the convex problem with the same objective value.
- Conversely, any feasible point  $(\mathbf{y}, t)$  in the convex problem leads to a feasible point  $\mathbf{x}$  in the original FP via  $\mathbf{x} = \mathbf{y}/t$ , also with the same objective value:

$$tf\left(\frac{\mathbf{y}}{t}\right) = \frac{f(\mathbf{x})}{g(\mathbf{x})}.$$

# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD**
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary

# Block-Coordinate Descent

## Block-Coordinate Descent (BCD) Method

- Also known as: Gauss-Seidel method, alternate minimization method.
- Objective: solve a difficult optimization problem by solving a sequence of simpler subproblems.

### Problem Formulation:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ & \text{subject to} && \mathbf{x}_i \in \mathcal{X}_i, \quad i = 1, \dots, n, \end{aligned}$$

where  $f$  is the (possibly nonconvex) objective function, each  $\mathcal{X}_i$  is a convex set, and the variable  $\mathbf{x}$  is partitioned into  $n$  blocks  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ .

### Method Description

- The iterative process produces a sequence of iterates  $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$  that converge to  $\mathbf{x}^*$ .
- The update rule optimizes the problem with respect to each block  $\mathbf{x}_i$  sequentially.
- At each outer iteration  $k$ , the method executes  $n$  inner iterations sequentially:

$$\mathbf{x}_i^{k+1} = \arg \min_{\mathbf{x}_i \in \mathcal{X}_i} f \left( \mathbf{x}_1^{k+1}, \dots, \mathbf{x}_{i-1}^{k+1}, \mathbf{x}_i, \mathbf{x}_{i+1}^k, \dots, \mathbf{x}_n^k \right), \quad i = 1, \dots, n.$$



# Algorithm

## Key points:

- Usefulness: Derive simple and practical algorithms.
- References: (Bertsekas 1999; Bertsekas and Tsitsiklis 1997; Beck 2017).

## BCD for separable problems

### Initialization:

- Choose initial point  $\mathbf{x}^0 = (\mathbf{x}_1^0, \dots, \mathbf{x}_n^0) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_n$ .
- Set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Execute  $n$  inner iterations sequentially:

$$\mathbf{x}_i^{k+1} = \arg \min_{\mathbf{x}_i \in \mathcal{X}_i} f \left( \mathbf{x}_1^{k+1}, \dots, \mathbf{x}_{i-1}^{k+1}, \mathbf{x}_i, \mathbf{x}_{i+1}^k, \dots, \mathbf{x}_n^k \right), \quad i = 1, \dots, n.$$

- 2  $k \leftarrow k + 1$

**Until:** convergence

## Monotonicity:

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k).$$

## Assumptions

- $f$  is continuously differentiable over the convex closed set  $\mathcal{X} = \mathcal{X}_1 \times \cdots \times \mathcal{X}_n$ .
- $f$  is blockwise strictly convex in each block variable  $\mathbf{x}_i$ .

## Convergence

Every limit point of the sequence  $\{\mathbf{x}^k\}$  is a stationary point of the original problem.

## References

(Bertsekas 1999; Bertsekas and Tsitsiklis 1997; Grippo and Sciandrone 2000).

# BCD: Parallel Updates

## Parallel Update (Jacobi Method)

- The method consists in executing  $n$  inner iterations in parallel instead of sequentially.
- The update rule is:

$$\mathbf{x}_i^{k+1} = \arg \min_{\mathbf{x}_i \in \mathcal{X}_i} f \left( \mathbf{x}_1^k, \dots, \mathbf{x}_{i-1}^k, \mathbf{x}_i, \mathbf{x}_{i+1}^k, \dots, \mathbf{x}_n^k \right), \quad i = 1, \dots, n.$$

## Jacobi Method

- The Jacobi method consists in a parallel update of the block variables.
- It is algorithmically attractive due to its potential for faster execution.

## Convergence Properties

- The Jacobi method does not enjoy nice convergence properties.
- Convergence is guaranteed if the mapping defined by  $T(\mathbf{x}) = \mathbf{x} - \gamma \nabla f(\mathbf{x})$  is a contraction for some  $\gamma$ .
- Reference: (Bertsekas 1999).

# BCD Example: Soft-Thresholding Operator

## Univariate Convex Optimization Problem

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{a}x - \mathbf{b}\|_2^2 + \lambda|x|$$

**Solution:**

$$x = \frac{1}{\|\mathbf{a}\|_2^2} \text{sign}(\mathbf{a}^\top \mathbf{b}) (|\mathbf{a}^\top \mathbf{b}| - \lambda)^+$$

- Sign function:

$$\text{sign}(u) = \begin{cases} +1 & u > 0 \\ 0 & u = 0 \\ -1 & u < 0 \end{cases}$$

- Positive part function:  $(\cdot)^+ = \max(0, \cdot)$

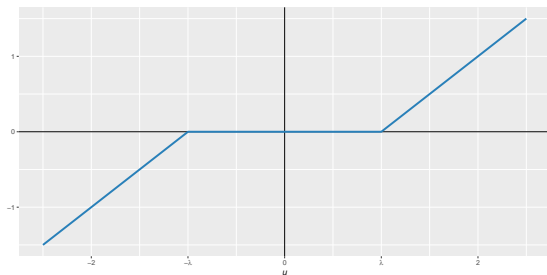
# BCD Example: Soft-Thresholding Operator

**Solution in compact form:**

$$\mathbf{x} = \frac{1}{\|\mathbf{a}\|_2^2} \mathcal{S}_\lambda \left( \mathbf{a}^\top \mathbf{b} \right)$$

**Soft-thresholding operator:**

$$\mathcal{S}_\lambda(u) = \text{sign}(u)(|u| - \lambda)^+$$



# BCD Example: $\ell_2 - \ell_1$ -Norm Minimization

## Problem Formulation:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

## Solution Approach

- Can be solved with a QP solver.
- Iterative algorithm via BCD with soft-thresholding operator (Zibulevsky and Elad 2010).

## BCD Method

- Variable partitioning: divide the variable into each constituent element  $\mathbf{x} = (x_1, \dots, x_n)$ .
- Sequence of problems at each iteration  $k = 0, 1, 2, \dots$  for each element  $i = 1, \dots, n$ :

$$\underset{x_i}{\text{minimize}} \quad \frac{1}{2} \left\| \mathbf{a}_i x_i - \tilde{\mathbf{b}}_i^k \right\|_2^2 + \lambda |x_i|$$

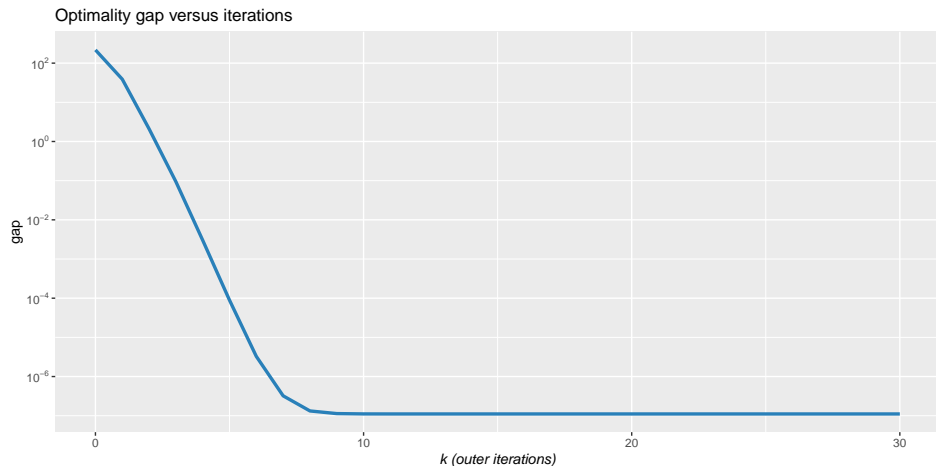
where  $\tilde{\mathbf{b}}_i^k \triangleq \mathbf{b} - \sum_{j < i} \mathbf{a}_j x_j^{k+1} - \sum_{j > i} \mathbf{a}_j x_j^k$ .

**Iterative Algorithm:** For  $k = 0, 1, 2, \dots$

$$x_i^{k+1} = \frac{1}{\|\mathbf{a}_i\|_2^2} \mathcal{S}_\lambda \left( \mathbf{a}_i^\top \tilde{\mathbf{b}}_i^k \right), \quad i = 1, \dots, n.$$

# BCD Example: $\ell_2 - \ell_1$ -Norm Minimization

Convergence of BCD for the  $\ell_2 - \ell_1$ -norm minimization:



# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM**
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary



# Majorization-Minimization (MM)

## Original Difficult Problem

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{X}\end{array}$$

where the objective function  $f$  and the feasible set  $\mathcal{X}$  are possibly nonconvex.

## Majorization-Minimization (MM) Method

- Approximates a difficult optimization problem by a sequence of simpler problems:

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} u(\mathbf{x}; \mathbf{x}^k) \quad k = 0, 1, 2, \dots$$

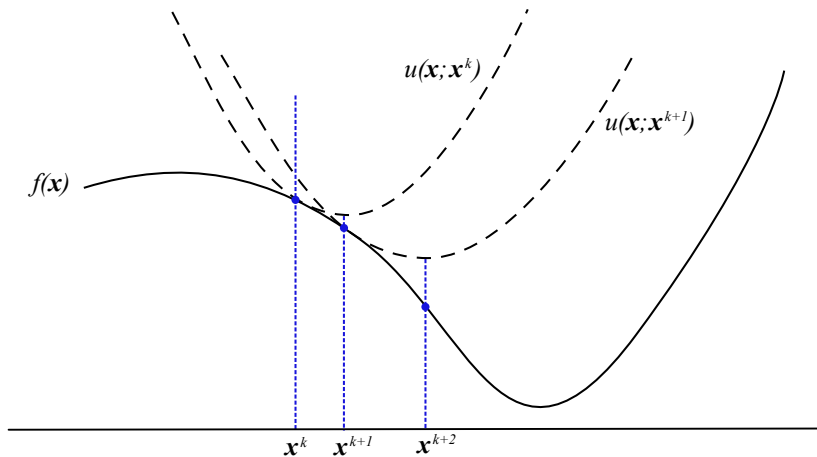
where  $u(\mathbf{x}; \mathbf{x}^k)$  is a surrogate function at iteration  $k$  that approximates  $f(\mathbf{x})$  around the current point  $\mathbf{x}^k$ .

- Produces a sequence of iterates  $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$  that converge to  $\mathbf{x}^*$ .

## References

- Concise tutorial: (Hunter and Lange 2004).
- Long tutorial with applications: (Sun, Babu, and Palomar 2017).
- Convergence analysis: (Razaviyayn, Hong, and Luo 2013).

Illustration of sequence of surrogate problems in MM:



## MM algorithm

### Initialization:

- Choose initial point  $\mathbf{x}^0 \in \mathcal{X}$ .
- Set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Construct majorizer of  $f(\mathbf{x})$  around current point  $\mathbf{x}^k$  as  $u(\mathbf{x}; \mathbf{x}^k)$ .
- 2 Obtain next iterate by solving the majorized problem:

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} u(\mathbf{x}; \mathbf{x}^k).$$

- 3  $k \leftarrow k + 1$

**Until:** convergence

## Conditions for the Surrogate Function:

- Upper-bound property:  $u(\mathbf{x}; \mathbf{x}^k) \geq f(\mathbf{x})$ .
- Touching property:  $u(\mathbf{x}^k; \mathbf{x}^k) = f(\mathbf{x}^k)$ .
- Tangent property:  $u(\mathbf{x}; \mathbf{x}^k)$  must be differentiable with  $\nabla u(\mathbf{x}; \mathbf{x}^k) = \nabla f(\mathbf{x})$ .

## Properties

- Monotonicity:  $f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k)$ .
- Convergence: If  $\mathcal{X}$  is convex, every limit point of the sequence  $\{\mathbf{x}^k\}$  is a stationary point of the original problem.
- If  $\mathcal{X}$  is nonconvex, then convergence must be studied on a case-by-case basis, e.g., (Song, Babu, and Palomar 2015; Sun, Babu, and Palomar 2017; Kumar et al. 2019, 2020).

## Majorizer Construction

- Objective: find an appropriate majorizer  $u(\mathbf{x}; \mathbf{x}^k)$  that satisfies the technical conditions and leads to a simpler surrogate problem.
- Techniques and examples: refer to (Sun, Babu, and Palomar 2017).

## MM Convergence Speed

- Issue: MM may require many iterations to converge if the surrogate function  $u(\mathbf{x}; \mathbf{x}^k)$  is not tight enough.
- Reason: Strict global upper-bound requirement.

## Acceleration Techniques

- Objective: Improve convergence speed.
- Popular technique: SQUAREM (Squared Iterative Methods for Accelerating EM-like Monotone Algorithms) (Varadhan and Roland 2008).

# MM Example: Nonnegative LS

## Problem Formulation

$$\underset{\mathbf{x} \geq \mathbf{0}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

where the parameters are

- $\mathbf{b} \in \mathbb{R}_+^m$  (nonnegative elements)
- $\mathbf{A} \in \mathbb{R}_{++}^{m \times n}$  (positive elements)

## Conventional LS Solution

Not applicable due to nonnegativity constraints:  $\mathbf{x}^* \neq (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ .

## Alternative Approach

- Use a QP solver: Standard method.
- Develop an iterative algorithm based on MM: More interesting approach.

# MM Example: Nonnegative LS

**Objective Function:**

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2.$$

**Majorizer**

$$u(\mathbf{x}; \mathbf{x}^k) = f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^k)^\top \Phi(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k)$$

- Gradient:  $\nabla f(\mathbf{x}^k) = \mathbf{A}^\top \mathbf{Ax}^k - \mathbf{A}^\top \mathbf{b}.$
- Matrix  $\Phi$ :  $\Phi(\mathbf{x}^k) = \text{Diag}\left(\frac{[\mathbf{A}^\top \mathbf{Ax}^k]_1}{x_1^k}, \dots, \frac{[\mathbf{A}^\top \mathbf{Ax}^k]_n}{x_n^k}\right).$

**Verification of Majorizer Properties**

- Upper-bound property:  $u(\mathbf{x}; \mathbf{x}^k) \geq f(\mathbf{x})$  (proved using Jensen's inequality)
- Touching property:  $u(\mathbf{x}^k; \mathbf{x}^k) = f(\mathbf{x}^k)$
- Tangent property:  $\nabla u(\mathbf{x}^k; \mathbf{x}^k) = \nabla f(\mathbf{x}^k)$

# MM Example: Nonnegative LS

## Sequence of Majorized Problems

$$\underset{\mathbf{x} \geq \mathbf{0}}{\text{minimize}} \quad \nabla f(\mathbf{x}^k)^\top \mathbf{x} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^k)^\top \Phi(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k)$$

with solution

$$\mathbf{x} = \mathbf{x}^k - \Phi(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k).$$

## Iterative Update:

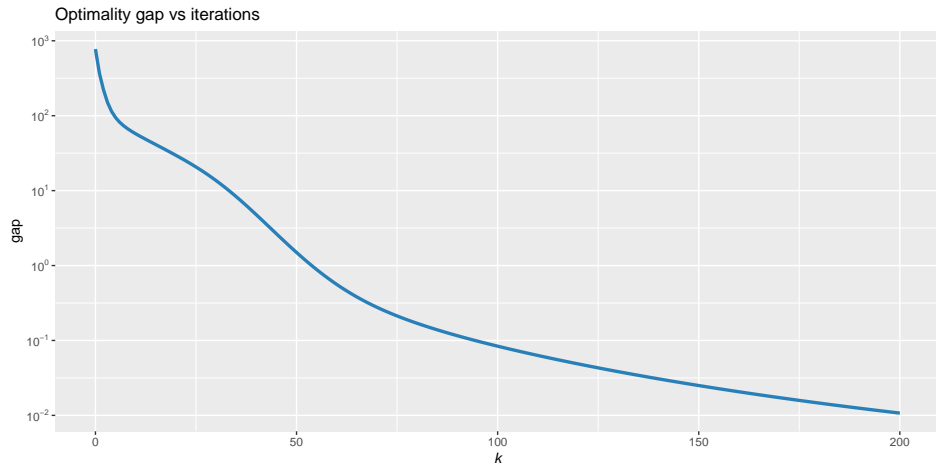
$$\mathbf{x}^{k+1} = \mathbf{c}^k \odot \mathbf{x}^k, \quad k = 0, 1, 2, \dots$$

where  $c_i^k = \frac{[\mathbf{A}^\top \mathbf{b}]_i}{[\mathbf{A}^\top \mathbf{A} \mathbf{x}^k]_i}$  and  $\odot$  denotes the elementwise product.



# MM Example: Nonnegative LS

Convergence of MM for nonnegative LS:



# MM Example: $\ell_2 - \ell_1$ -Norm Minimization

$\ell_2 - \ell_1$ -Norm Minimization Problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

## Solution Methods

- Can be solved via BCD or a QP solver.
- We will develop an iterative algorithm based on MM.

## MM for $\ell_2 - \ell_1$ -Norm Minimization

- The variable  $\mathbf{x}$  is partitioned into elements  $(x_1, \dots, x_n)$ .
- The majorizer function is:

$$u(\mathbf{x}; \mathbf{x}^k) = \frac{\kappa}{2} \|\mathbf{x} - \bar{\mathbf{x}}^k\|_2^2 + \lambda \|\mathbf{x}\|_1 + \text{constant},$$

$$\text{where } \bar{\mathbf{x}}^k = \mathbf{x}^k - \frac{1}{\kappa} \mathbf{A}^T (\mathbf{Ax}^k - \mathbf{b}).$$

- Verification as a majorizer is given by:

$$u(\mathbf{x}; \mathbf{x}^k) = f(\mathbf{x}) + \text{dist}(\mathbf{x}, \mathbf{x}^k),$$

$$\text{where } \text{dist}(\mathbf{x}, \mathbf{x}^k) = \frac{\kappa}{2} \|\mathbf{x} - \mathbf{x}^k\|_2^2 - \frac{1}{2} \|\mathbf{Ax} - \mathbf{Ax}^k\|_2^2 \text{ and } \kappa > \lambda_{\max}(\mathbf{A}^T \mathbf{A}).$$

## MM Example: $\ell_2 - \ell_1$ -Norm Minimization

**Sequence of Majorized Problems:** For  $k = 0, 1, 2, \dots$

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{\kappa}{2} \|\mathbf{x} - \bar{\mathbf{x}}^k\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

**MM Iterative Algorithm:**

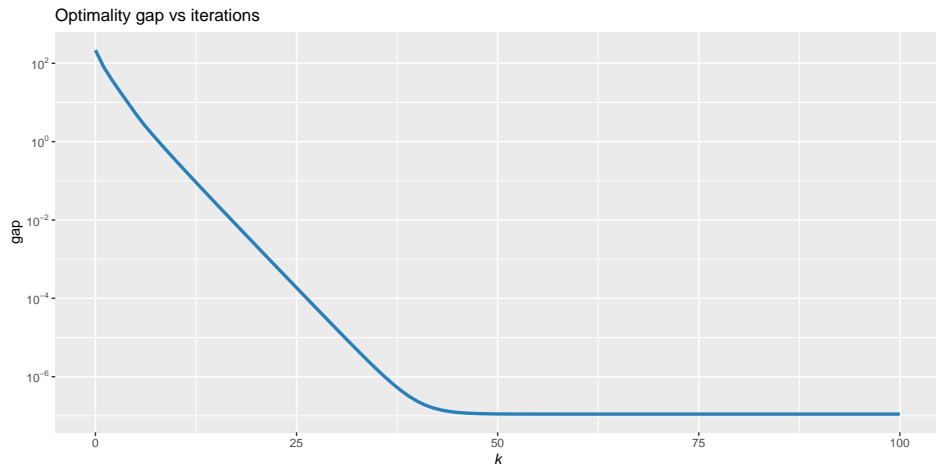
$$\mathbf{x}^{k+1} = \mathcal{S}_{\lambda/\kappa}(\bar{\mathbf{x}}^k), \quad k = 0, 1, 2, \dots$$

where  $\mathcal{S}_{\lambda/\kappa}(\cdot)$  is the soft-thresholding operator:

$$\mathcal{S}_{\lambda/\kappa}(z) = \text{sign}(z) \max(|z| - \lambda/\kappa, 0).$$

# MM Example: $\ell_2 - \ell_1$ -Norm Minimization

Convergence of MM for  $\ell_2 - \ell_1$ -norm minimization:



# Block MM: Combining BCD and MM

## Objective

- Address situations where both the original problem and direct application of MM are too difficult to solve.

## Approach

- Combine Block-Coordinate Descent (BCD) and Majorization-Minimization (MM).

## Original Problem

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ \text{subject to} & \mathbf{x}_i \in \mathcal{X}_i, \quad i = 1, \dots, n\end{array}$$

- Partitioning: Variables are partitioned into  $n$  blocks  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ .
- Constraints: Each block  $\mathbf{x}_i$  is separately constrained.

## Idea

- Solve the problem block by block as in BCD, but majorize each block  $f(\mathbf{x}_i)$  with a surrogate function  $u(\mathbf{x}_i; \mathbf{x}^k)$ .

# Block MM: Combining BCD and MM

## Procedure

**Initialization:** Start with an initial guess  $\mathbf{x}^0 = (\mathbf{x}_1^0, \dots, \mathbf{x}_n^0)$ .

**Iterative process:** For each outer iteration  $k = 0, 1, 2, \dots$ :

- For each block  $i = 1, \dots, n$ :
  - Majorize: Construct a surrogate function  $u(\mathbf{x}_i; \mathbf{x}^k)$  for the block  $f(\mathbf{x}_i)$ .
  - Update: Solve the majorized problem for the block:

$$\mathbf{x}_i^{k+1} = \arg \min_{\mathbf{x}_i \in \mathcal{X}_i} u(\mathbf{x}_i; \mathbf{x}^k).$$

- Update the full variable:  $\mathbf{x}^{k+1} = (\mathbf{x}_1^{k+1}, \dots, \mathbf{x}_n^{k+1})$ .

## References

(Razaviyayn, Hong, and Luo 2013) (Sun, Babu, and Palomar 2017).

# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA**
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary

# Successive Convex Approximation (SCA)

## Original Difficult Problem:

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{X}\end{array}$$

where  $f$  is the (possibly nonconvex) objective function and  $\mathcal{X}$  is the convex feasible set.

## Successive Convex Approximation (SCA) Method (Scutari et al. 2014):

- Approximates a difficult optimization problem by a sequence of simpler problems:

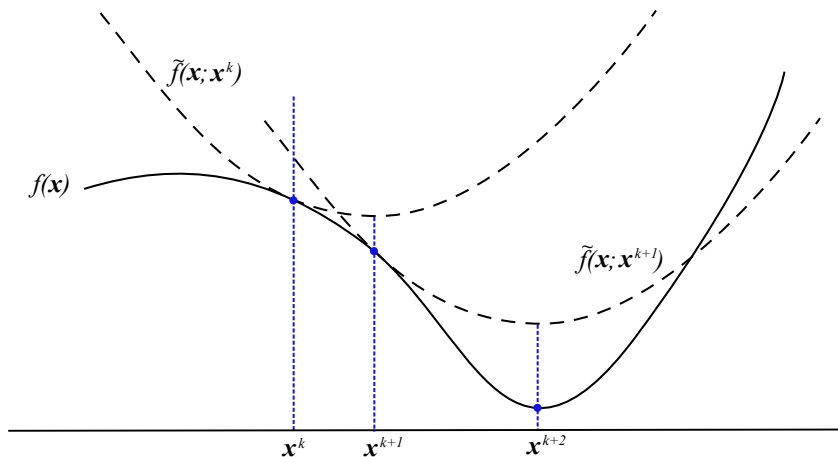
$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \tilde{f}(\mathbf{x}; \mathbf{x}^k), \quad k = 0, 1, 2, \dots$$

- The function  $\tilde{f}(\mathbf{x}; \mathbf{x}^k)$  is a surrogate function at iteration  $k$  that approximates  $f(\mathbf{x})$  around the current point  $\mathbf{x}^k$ .
- The method produces a sequence of iterates  $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$  that converge to  $\mathbf{x}^*$ .
- For convergence, a smoothing step is needed to avoid oscillations ( $\gamma^k \in (0, 1]$ ):

$$\begin{aligned}\hat{\mathbf{x}}^{k+1} &= \arg \min_{\mathbf{x} \in \mathcal{X}} \tilde{f}(\mathbf{x}; \mathbf{x}^k) \\ \mathbf{x}^{k+1} &= \mathbf{x}^k + \gamma^k (\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k)\end{aligned} \quad k = 0, 1, 2, \dots$$



Illustration of sequence of surrogate problems in SCA:



# SCA: Convergence

## Conditions for the Surrogate Function $\tilde{f}(\mathbf{x}; \mathbf{x}^k)$

- Must be strongly convex on the feasible set  $\mathcal{X}$ .
- Must be differentiable with  $\nabla \tilde{f}(\mathbf{x}; \mathbf{x}^k) = \nabla f(\mathbf{x})$ .

## Stepsize Rules for $\{\gamma^k\}$

- *Bounded stepsize*:  $\gamma^k$  values are sufficiently small (difficult to use in practice).
- *Backtracking line search*: Effective in terms of iterations but costly.
- *Diminishing stepsize*: satisfying  $\sum_{k=1}^{\infty} \gamma^k = +\infty$  and  $\sum_{k=1}^{\infty} (\gamma^k)^2 < +\infty$ .
  - Example 1:  $\gamma^{k+1} = \gamma^k (1 - \epsilon \gamma^k)$ ,  $\gamma^0 < 1/\epsilon$ ,  $\epsilon \in (0, 1)$ .
  - Example 2:  $\gamma^{k+1} = \frac{\gamma^k + \alpha^k}{1 + \beta^k}$ ,  $\gamma^0 = 1$ ,  $\alpha^k$  and  $\beta^k$  satisfy  $0 \leq \alpha^k \leq \beta^k$  and  $\alpha^k / \beta^k \rightarrow 0$ .

## Examples of $\alpha^k$ and $\beta^k$

- $\alpha^k = \alpha$  or  $\alpha^k = \log(k)^\alpha$ .
- $\beta^k = \beta k$  or  $\beta^k = \beta \sqrt{k}$ .
- Constants  $\alpha \in (0, 1)$ ,  $\beta \in (0, 1)$ , and  $\alpha \leq \beta$ .

## Advantages of SCA

- Surrogate function is convex by construction.
- Easier to construct a convex surrogate function compared to MM.

## SCA algorithm

### Initialization:

- Choose initial point  $\mathbf{x}^0 \in \mathcal{X}$ , sequence  $\{\gamma^k\}$ , and set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Construct surrogate of  $f(\mathbf{x})$  around current point  $\mathbf{x}^k$  as  $\tilde{f}(\mathbf{x}; \mathbf{x}^k)$ .
- 2 Obtain intermediate point by solving the surrogate convex problem:

$$\hat{\mathbf{x}}^{k+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \tilde{f}(\mathbf{x}; \mathbf{x}^k).$$

- 3 Obtain next iterate by averaging the intermediate point with the previous one:

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \gamma^k (\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k).$$

- 4  $k \leftarrow k + 1$

**Until:** convergence

# Gradient Descent and Newton's Method as SCA

## Gradient Descent Method as SCA

- Unconstrained problem

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}),$$

- One possible SCA surrogate function is

$$\tilde{f}(\mathbf{x}; \mathbf{x}^k) = f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2\alpha^k} \|\mathbf{x} - \mathbf{x}^k\|^2.$$

- To minimize the surrogate function, set its gradient to zero:

$$\nabla \tilde{f}(\mathbf{x}; \mathbf{x}^k) = \nabla f(\mathbf{x}^k) + \frac{1}{\alpha^k} (\mathbf{x} - \mathbf{x}^k) = 0.$$

- Solving for  $\mathbf{x}$  yields:

$$\mathbf{x} = \mathbf{x}^k - \alpha^k \nabla f(\mathbf{x}^k).$$

- The resulting iteration process coincides with the gradient descent method:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha^k \nabla f(\mathbf{x}^k), \quad k = 0, 1, 2, \dots$$

## Newton's Method as SCA

- By including second-order information, the surrogate function becomes:

$$\tilde{f}(\mathbf{x}; \mathbf{x}^k) = f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2\alpha^k} (\mathbf{x} - \mathbf{x}^k)^\top \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k).$$

- To minimize the surrogate function, set its gradient to zero:

$$\nabla \tilde{f}(\mathbf{x}; \mathbf{x}^k) = \nabla f(\mathbf{x}^k) + \frac{1}{\alpha^k} \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k) = 0.$$

- Solving for  $\mathbf{x}$  yields:

$$\mathbf{x} = \mathbf{x}^k - \alpha^k \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k).$$

- The resulting iteration process coincides with Newton's method:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha^k \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k), \quad k = 0, 1, 2, \dots$$

## Partitioned Variables in SCA

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ \text{subject to} & \mathbf{x}_i \in \mathcal{X}_i, \quad i = 1, \dots, n,\end{array}$$

with variables are partitioned into  $n$  separate blocks:  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ .

## Parallel Updates in SCA

- Unlike BCD or MM, SCA updates variables in parallel with surrogates  $\tilde{f}_i(\mathbf{x}_i; \mathbf{x}^k)$ .
- The update process for each block  $i$  is:

$$\begin{aligned}\hat{\mathbf{x}}_i^{k+1} &= \arg \min_{\mathbf{x}_i \in \mathcal{X}_i} \tilde{f}_i(\mathbf{x}_i; \mathbf{x}^k) \\ \mathbf{x}_i^{k+1} &= \mathbf{x}_i^k + \gamma^k (\hat{\mathbf{x}}_i^{k+1} - \mathbf{x}_i^k)\end{aligned} \quad i = 1, \dots, n, \quad k = 0, 1, 2, \dots$$

where  $\{\gamma^k\}$  is a properly designed sequence with  $\gamma^k \in (0, 1]$ .

## Advantages of Parallel Updates

- Efficiently handles large-scale problems by updating multiple variables simultaneously.
- Reduces computational time compared to sequential updates in BCD or block MM.

## Technical Conditions for the Surrogate Function (Scutari et al. 2014)

- Must be strongly convex on the feasible set  $\mathcal{X}$ .
- Must be differentiable with  $\nabla \tilde{f}(\mathbf{x}; \mathbf{x}^k) = \nabla f(\mathbf{x})$ .

## Stepsize Rules for $\{\gamma^k\}$

- *Bounded stepsize*:  $\gamma^k$  values are sufficiently small.
- *Backtracking line search*: effective but requires multiple evaluations per iteration.
- *Diminishing stepsize*: satisfying  $\sum_{k=1}^{\infty} \gamma^k = +\infty$  and  $\sum_{k=1}^{\infty} (\gamma^k)^2 < +\infty$ .

## Theoretical Convergence

- SCA enjoys strong theoretical convergence properties.
- Convergence results are detailed in (Scutari et al. 2014).

## Convergence of SCA

- Suppose the surrogate function  $\tilde{f}(\mathbf{x}; \mathbf{x}^k)$  (or each  $\tilde{f}_i(\mathbf{x}_i; \mathbf{x}^k)$  in the parallel version) satisfies the required technical conditions.
- If  $\{\gamma^k\}$  is chosen according to the bounded stepsize, diminishing rule, or backtracking line search, then the sequence  $\{\mathbf{x}^k\}$  converges to a stationary point of the original problem.

# SCA Example: $\ell_2 - \ell_1$ -Norm Minimization

$\ell_2 - \ell_1$ -Norm Minimization Problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

## Solution Methods

- Can be solved via BCD, MM, or a QP solver.
- We will develop an iterative algorithm based on SCA.

## Parallel SCA for $\ell_2 - \ell_1$ -Norm Minimization

- The variable  $\mathbf{x}$  is partitioned into elements  $(x_1, \dots, x_n)$ .
- Surrogate functions:

$$\tilde{f}(\mathbf{x}_i; \mathbf{x}^k) = \frac{1}{2} \|\mathbf{a}_i x_i - \tilde{\mathbf{b}}_i^k\|_2^2 + \lambda |x_i| + \frac{\tau}{2} (x_i - x_i^k)^2,$$

where  $\tilde{\mathbf{b}}_i^k = \mathbf{b} - \sum_{j \neq i} \mathbf{a}_j x_j^k$ .

**Sequence of surrogate problems:** For  $k = 0, 1, 2, \dots$  and  $i = 1, \dots, n$ :

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{a}_i x_i - \tilde{\mathbf{b}}_i^k\|_2^2 + \lambda |x_i| + \tau (x_i - x_i^k)^2$$



## SCA Example: $\ell_2 - \ell_1$ -Norm Minimization

**SCA iterative algorithm:**

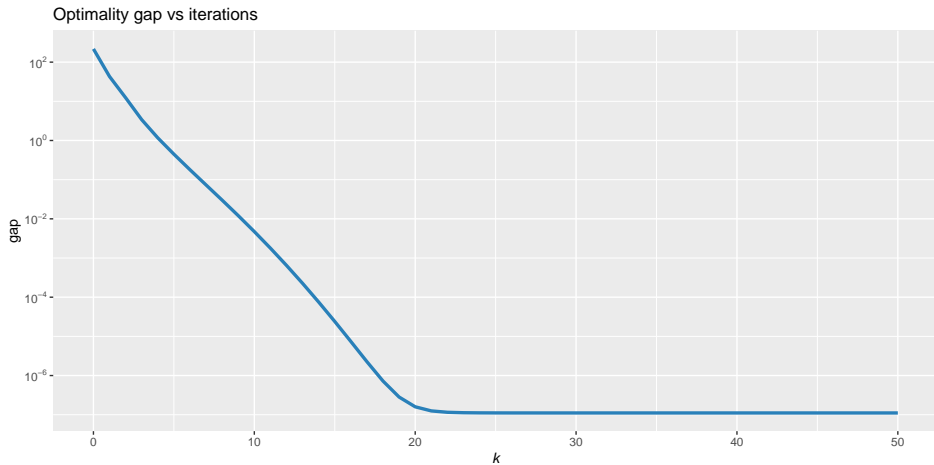
$$\begin{aligned}\hat{x}_i^{k+1} &= \frac{1}{\tau + \|\mathbf{a}_i\|^2} \mathcal{S}_\lambda \left( \mathbf{a}_i^\top \tilde{\mathbf{b}}_i^k + \tau x_i^k \right) \\ x_i^{k+1} &= x_i^k + \gamma^k \left( \hat{x}_i^{k+1} - x_i^k \right)\end{aligned} \quad i = 1, \dots, n, \quad k = 0, 1, 2, \dots$$

where  $\mathcal{S}_\lambda(\cdot)$  is the soft-thresholding operator:

$$\mathcal{S}_\lambda(z) = \text{sign}(z) \max(|z| - \lambda, 0).$$

# SCA Example: $\ell_2 - \ell_1$ -Norm Minimization

Convergence of SCA for the  $\ell_2 - \ell_1$ -norm minimization:



# SCA Example: Dictionary Learning

## Dictionary Learning Problem

$$\begin{aligned} & \underset{\mathbf{D}, \mathbf{X}}{\text{minimize}} && \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1 \\ & \text{subject to} && \|[\mathbf{D}]_{:,i}\| \leq 1, \quad i = 1, \dots, m. \end{aligned}$$

- $\|\mathbf{D}\|_F$ : frobenius norm of  $\mathbf{D}$
- $\|\mathbf{X}\|_1$ : elementwise  $\ell_1$ -norm of  $\mathbf{X}$

## Matrix Definitions

- $\mathbf{D}$ : dictionary matrix (fat matrix with columns explaining the columns of  $\mathbf{Y}$ ).
- $\mathbf{X}$ : sparse matrix selecting a few columns of the dictionary.

## Bi-Convex Nature

- The problem is not jointly convex in  $(\mathbf{D}, \mathbf{X})$ , but it is bi-convex.
- For fixed  $\mathbf{D}$ , the problem is convex in  $\mathbf{X}$ .
- For fixed  $\mathbf{X}$ , the problem is convex in  $\mathbf{D}$ .

## Solution Methods

- BCD: updates  $\mathbf{D}$  and  $\mathbf{X}$  sequentially.
- SCA: allows parallel updates of  $\mathbf{D}$  and  $\mathbf{X}$ .

# SCA Example: Dictionary Learning

**SCA Approach:** Surrogate functions:

$$\begin{aligned}\tilde{f}_1(\mathbf{D}; \mathbf{X}^k) &= \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{X}^k\|_F^2 \\ \tilde{f}_2(\mathbf{X}; \mathbf{D}^k) &= \frac{1}{2} \|\mathbf{Y} - \mathbf{D}^k\mathbf{X}\|_F^2\end{aligned}$$

## Resulting Convex Problems

- Normalized least squares (LS) problem for  $\mathbf{D}$ :

$$\begin{aligned}\underset{\mathbf{D}}{\text{minimize}} \quad & \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{X}^k\|_F^2 \\ \text{subject to} \quad & \|[\mathbf{D}]_{:,i}\| \leq 1, \quad i = 1, \dots, m\end{aligned}$$

- Matrix version of the  $\ell_2 - \ell_1$ -norm problem for  $\mathbf{X}$ :

$$\underset{\mathbf{X}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}^k\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1$$

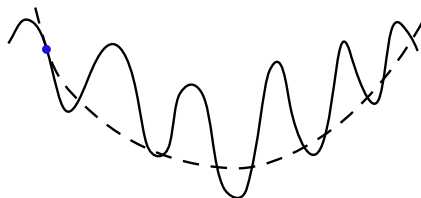
which can be further decomposed into a set of vectorized  $\ell_2 - \ell_1$ -norm problems for each column of  $\mathbf{X}$ .

## Surrogate Function

- MM (Majorization-Minimization):
  - Requires the surrogate function to be a global upper bound.
  - The surrogate function need not be convex.
  - Can be difficult to derive and too restrictive in some cases.
- SCA (Successive Convex Approximation):
  - Relaxes the upper-bound condition.
  - Requires the surrogate function to be strongly convex.



MM



SCA

**Constraint Set:** In principle, both require the feasible set  $\mathcal{X}$  to be convex.

- MM:
  - Convergence can be extended to nonconvex  $\mathcal{X}$  on a case-by-case basis.
  - Examples of nonconvex  $\mathcal{X}$  handled by MM: (Song, Babu, and Palomar 2015; Sun, Babu, and Palomar 2017; Kumar et al. 2019, 2020).
- SCA:
  - Cannot directly handle nonconvex  $\mathcal{X}$  (Scutari et al. 2014).
  - Some extensions allow for successive convexification of  $\mathcal{X}$ , but at the expense of a more complex algorithm (Scutari and Sun 2018).

**Schedule of Updates:** Both can handle separable variables  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ .

- MM:
  - Requires a sequential update for block variables (Razaviyayn, Hong, and Luo 2013; Sun, Babu, and Palomar 2017).
- SCA:
  - Naturally implements a parallel update, which is more amenable for distributed implementations.

# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM**
- 9 Numerical Comparison
- 10 Summary

## Alternating Direction Method of Multipliers (ADMM)

- Practical algorithm resembling BCD but can handle coupled block variables in constraints.
- Detailed in (S. Boyd et al. 2010) and (Beck 2017).

## Convex Optimization Problem

$$\begin{array}{ll}\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} & f(\mathbf{x}) + g(\mathbf{z}) \\ \text{subject to} & \mathbf{Ax} + \mathbf{Bz} = \mathbf{c},\end{array}$$

Observe that the variables  $\mathbf{x}$  and  $\mathbf{z}$  are coupled via the constraint  $\mathbf{Ax} + \mathbf{Bz} = \mathbf{c}$ .



# First Attempt: Dual Ascent Method

First attempt to decouple the variables.

## Dual Ascent Method

- Updates dual variable  $\mathbf{y}$  via gradient method.
- Solves Lagrangian for given  $\mathbf{y}$ :

$$\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} \quad L(\mathbf{x}, \mathbf{z}; \mathbf{y}) \triangleq f(\mathbf{x}) + g(\mathbf{z}) + \mathbf{y}^T (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c})$$

- Decouples into two separate problems over  $\mathbf{x}$  and  $\mathbf{z}$ :

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{x}} f(\mathbf{x}) + (\mathbf{y}^k)^T \mathbf{A}\mathbf{x}$$

$$\mathbf{z}^{k+1} = \arg \min_{\mathbf{z}} g(\mathbf{z}) + (\mathbf{y}^k)^T \mathbf{B}\mathbf{z} \quad k = 0, 1, 2, \dots$$

$$\mathbf{y}^{k+1} = \mathbf{y}^k + \alpha^k (\mathbf{A}\mathbf{x}^{k+1} + \mathbf{B}\mathbf{z}^{k+1} - \mathbf{c})$$

- Requires many technical assumptions and is often slow.

## Second Attempt: Method of Multipliers

Second attempt to decouple the variables.

### Method of Multipliers

- Uses the augmented Lagrangian:

$$L_{\rho}(\mathbf{x}, \mathbf{z}; \mathbf{y}) \triangleq f(\mathbf{x}) + g(\mathbf{z}) + \mathbf{y}^{\top} (\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}\|_2^2$$

- The algorithm is:

$$\begin{aligned} (\mathbf{x}^{k+1}, \mathbf{z}^{k+1}) &= \arg \min_{\mathbf{x}, \mathbf{z}} L_{\rho}(\mathbf{x}, \mathbf{z}; \mathbf{y}^k) \\ \mathbf{y}^{k+1} &= \mathbf{y}^k + \rho (\mathbf{Ax}^{k+1} + \mathbf{Bz}^{k+1} - \mathbf{c}) \end{aligned} \quad k = 0, 1, 2, \dots$$

- It converges under more relaxed conditions but cannot decouple  $\mathbf{x}$  and  $\mathbf{z}$  due to the  $\|\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}\|_2^2$  term.

## Third Attempt: ADMM

Third and final attempt to decouple the variables.

### Alternating Direction Method of Multipliers (ADMM)

- Combines features of dual decomposition and the method of multipliers.
- Minimizes the augmented Lagrangian with a BCD-like method:

$$\begin{aligned}\mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} L_{\rho}(\mathbf{x}, \mathbf{z}^k; \mathbf{y}^k) \\ \mathbf{z}^{k+1} &= \arg \min_{\mathbf{z}} L_{\rho}(\mathbf{x}^{k+1}, \mathbf{z}; \mathbf{y}^k) \\ \mathbf{y}^{k+1} &= \mathbf{y}^k + \rho (\mathbf{A}\mathbf{x}^{k+1} + \mathbf{B}\mathbf{z}^{k+1} - \mathbf{c})\end{aligned} \quad k = 0, 1, 2, \dots$$

- Successfully decouples the primal variables  $\mathbf{x}$  and  $\mathbf{z}$ .
- Achieves faster convergence with fewer technical conditions.

### Scaled Dual Variable Form

It is common to express the ADMM updates using the scaled dual variable  $\mathbf{u}^k = \mathbf{y}^k / \rho$ .

## ADMM algorithm

### Initialization:

- Choose initial point  $(\mathbf{x}^0, \mathbf{z}^0)$ ,  $\rho$ , and set  $k \leftarrow 0$ .

### Repeat ( $k$ th iteration):

- 1 Iterate primal and dual variables:

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{x}} f(\mathbf{x}) + \frac{\rho}{2} \left\| \mathbf{Ax} + \mathbf{Bz}^k - \mathbf{c} + \mathbf{u}^k \right\|_2^2$$

$$\mathbf{z}^{k+1} = \arg \min_{\mathbf{z}} g(\mathbf{z}) + \frac{\rho}{2} \left\| \mathbf{Ax}^{k+1} + \mathbf{Bz} - \mathbf{c} + \mathbf{u}^k \right\|_2^2$$

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \left( \mathbf{Ax}^{k+1} + \mathbf{Bz}^{k+1} - \mathbf{c} \right);$$

- 2  $k \leftarrow k + 1$

**Until:** convergence

# ADMM: Convergence\*

## Assumptions

- $f(\mathbf{x})$  and  $g(\mathbf{z})$  are convex.
- Both the  $\mathbf{x}$ -update and the  $\mathbf{z}$ -update are solvable.
- The Lagrangian has a saddle point.

## Convergence of ADMM

- Residual convergence (feasibility convergence of iterates):  $\mathbf{Ax}^k + \mathbf{Bz}^k - \mathbf{c} \rightarrow \mathbf{0}$  as  $k \rightarrow \infty$ .
- Objective convergence:  $f(\mathbf{x}) + g(\mathbf{z}) \rightarrow p^*$  as  $k \rightarrow \infty$ .
- Dual variable convergence:  $\mathbf{y}^k \rightarrow \mathbf{y}^*$  as  $k \rightarrow \infty$ .
- Detailed analysis in (S. Boyd et al. 2010) and references therein.

## Practical Considerations

- The sequences  $\{\mathbf{x}^k\}$  and  $\{\mathbf{z}^k\}$  need not converge to optimal values without additional assumptions.
- ADMM can be slow to converge to high accuracy.
- It often converges to modest accuracy within a few tens of iterations, which is sufficient for many practical applications.

# ADMM Example: Constrained Convex Optimization

## Generic Convex Optimization Problem:

$$\begin{array}{ll}\underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{X},\end{array}$$

where  $f$  is convex and  $\mathcal{X}$  is a convex set.

## Using ADMM to Transform the Problem

- Define  $g$  as the indicator function of the feasible set  $\mathcal{X}$ :

$$g(\mathbf{x}) \triangleq \begin{cases} 0 & \mathbf{x} \in \mathcal{X} \\ +\infty & \text{otherwise,} \end{cases}$$

- Formulate the equivalent problem:

$$\begin{array}{ll}\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} & f(\mathbf{x}) + g(\mathbf{z}) \\ \text{subject to} & \mathbf{x} - \mathbf{z} = \mathbf{0}.\end{array}$$

# ADMM Example: Constrained Convex Optimization

**ADMM Algorithm for the Transformed Problem:** The update rules are:

$$\begin{aligned}\mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} f(\mathbf{x}) + \frac{\rho}{2} \left\| \mathbf{x} - \mathbf{z}^k + \mathbf{u}^k \right\|_2^2 \\ \mathbf{z}^{k+1} &= \left[ \mathbf{x}^{k+1} + \mathbf{u}^k \right]_{\mathcal{X}} \\ \mathbf{u}^{k+1} &= \mathbf{u}^k + \left( \mathbf{x}^{k+1} - \mathbf{z}^{k+1} \right)\end{aligned} \quad k = 0, 1, 2, \dots$$

where  $[\cdot]_{\mathcal{X}}$  denotes projection on the set  $\mathcal{X}$ .

## Explanation of Steps

- The  $\mathbf{x}$ -update minimizes  $f(\mathbf{x})$  with a quadratic penalty term.
- The  $\mathbf{z}$ -update projects  $\mathbf{x}^{k+1} + \mathbf{u}^k$  onto the set  $\mathcal{X}$ .
- The  $\mathbf{u}$ -update updates the scaled dual variable  $\mathbf{u}$ .

## Benefits of this Approach

- Transforms a constrained optimization problem into an unconstrained one.
- Leverages the efficiency of ADMM for solving the problem.
- Allows for the use of projection operations to handle constraints.

# ADMM Example: $\ell_2 - \ell_1$ -Norm Minimization

$\ell_2 - \ell_1$ -Norm Minimization Problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

Reformulated Problem for ADMM:

$$\begin{aligned} &\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} && \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{z}\|_1 \\ &\text{subject to} && \mathbf{x} - \mathbf{z} = \mathbf{0}. \end{aligned}$$



# ADMM Example: $\ell_2 - \ell_1$ -Norm Minimization

## ADMM Algorithm

- The  $\mathbf{x}$ -update solves the unconstrained QP:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \frac{\rho}{2} \|\mathbf{x} - \mathbf{z} + \mathbf{u}\|_2^2,$$

which has the solution:

$$\mathbf{x} = \left( \mathbf{A}^\top \mathbf{A} + \rho \mathbf{I} \right)^{-1} \left( \mathbf{A}^\top \mathbf{b} + \rho(\mathbf{z} - \mathbf{u}) \right).$$

- The  $\mathbf{z}$ -update solves:

$$\underset{\mathbf{z}}{\text{minimize}} \quad \frac{\rho}{2} \|\mathbf{x} - \mathbf{z} + \mathbf{u}\|_2^2 + \lambda \|\mathbf{z}\|_1,$$

which has the solution using the soft-thresholding operator  $\mathcal{S}_{\lambda/\rho}(\cdot)$ :

$$\mathbf{z} = \mathcal{S}_{\lambda/\rho}(\mathbf{x} + \mathbf{u}).$$

- The  $\mathbf{u}$ -update for the scaled dual variable is:

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \left( \mathbf{x}^{k+1} - \mathbf{z}^{k+1} \right).$$

# ADMM Example: $\ell_2 - \ell_1$ -Norm Minimization

## ADMM Iterative Algorithm

The update rules are:

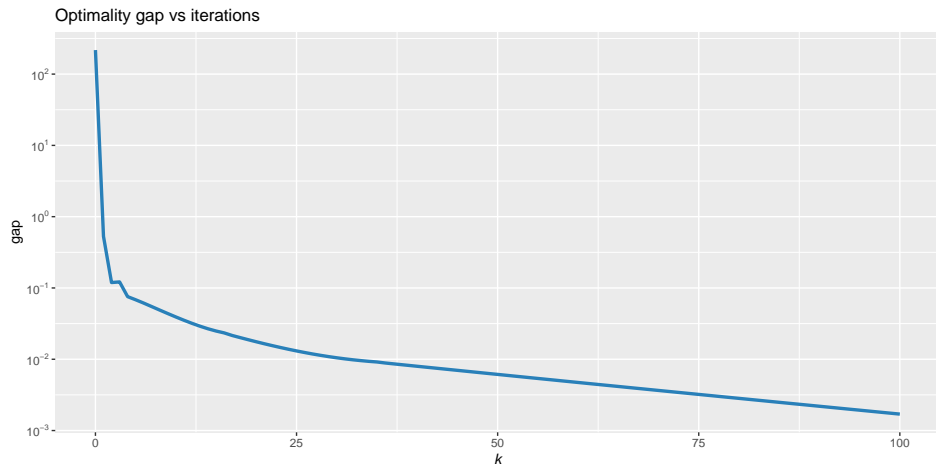
$$\begin{aligned}\mathbf{x}^{k+1} &= \left(\mathbf{A}^\top \mathbf{A} + \rho \mathbf{I}\right)^{-1} \left(\mathbf{A}^\top \mathbf{b} + \rho \left(\mathbf{z}^k - \mathbf{u}^k\right)\right) \\ \mathbf{z}^{k+1} &= \mathcal{S}_{\lambda/\rho} \left(\mathbf{x}^{k+1} + \mathbf{u}^k\right) \\ \mathbf{u}^{k+1} &= \mathbf{u}^k + \left(\mathbf{x}^{k+1} - \mathbf{z}^{k+1}\right)\end{aligned}\quad k = 0, 1, 2, \dots$$

where  $\mathcal{S}_{\lambda/\rho}(z)$  is the soft-thresholding operator:

$$\mathcal{S}_{\lambda/\rho}(z) = \text{sign}(z) \max(|z| - \lambda/\rho, 0).$$

# ADMM Example: $\ell_2 - \ell_1$ -Norm Minimization

Convergence of ADMM for the  $\ell_2 - \ell_1$ -norm minimization:



# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison**
- 10 Summary

# Numerical Comparison

$\ell_2 - \ell_1$ -Norm Minimization Problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

Iterates of Various Algorithms:

- BCD (Gauss-Seidel) iterates:

$$\mathbf{x}^{k+1} = \mathcal{S}_{\frac{\lambda}{\text{diag}(\mathbf{A}^T \mathbf{A})}} \left( \mathbf{x}^k - \frac{\mathbf{A}^T (\mathbf{Ax}^{(k,i)} - \mathbf{b})}{\text{diag}(\mathbf{A}^T \mathbf{A})} \right), \quad i = 1, \dots, n, \quad k = 0, 1, 2, \dots$$

where  $\mathbf{x}^{(k,i)} \triangleq (x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i^k, \dots, x_n^k)$

- Parallel BCD (Jacobi) iterates:

$$\mathbf{x}^{k+1} = \mathcal{S}_{\frac{\lambda}{\text{diag}(\mathbf{A}^T \mathbf{A})}} \left( \mathbf{x}^k - \frac{\mathbf{A}^T (\mathbf{Ax}^k - \mathbf{b})}{\text{diag}(\mathbf{A}^T \mathbf{A})} \right), \quad i = 1, \dots, n, \quad k = 0, 1, 2, \dots$$

## Iterates of Various Algorithms (cont'd):

- MM iterates:

$$\mathbf{x}^{k+1} = \mathcal{S}_{\frac{\lambda}{\kappa}} \left( \mathbf{x}^k - \frac{1}{\kappa} \mathbf{A}^T (\mathbf{A} \mathbf{x}^k - \mathbf{b}) \right), \quad k = 0, 1, 2, \dots$$

- Accelerated MM iterates:

$$\mathbf{r}^k = R(\mathbf{x}^k) \triangleq \text{MM}(\mathbf{x}^k) - \mathbf{x}^k$$

$$\mathbf{v}^k = R(\text{MM}(\mathbf{x}^k)) - R(\mathbf{x}^k)$$

$$\alpha^k = -\max \left( 1, \|\mathbf{r}^k\|_2 / \|\mathbf{v}^k\|_2 \right) \quad k = 0, 1, 2, \dots$$

$$\mathbf{y}^k = \mathbf{x}^k - \alpha^k \mathbf{r}^k$$

$$\mathbf{x}^{k+1} = \text{MM}(\mathbf{y}^k)$$

## Iterates of Various Algorithms (cont'd):

- **SCA iterates:**

$$\begin{aligned}\hat{\mathbf{x}}^{k+1} &= \mathcal{S}_{\frac{\lambda}{\tau + \text{diag}(\mathbf{A}^\top \mathbf{A})}} \left( \mathbf{x}^k - \frac{\mathbf{A}^\top (\mathbf{A} \mathbf{x}^k - \mathbf{b})}{\tau + \text{diag}(\mathbf{A}^\top \mathbf{A})} \right) & k = 0, 1, 2, \dots \\ \mathbf{x}^{k+1} &= \gamma^k \hat{\mathbf{x}}^{k+1} + (1 - \gamma^k) \mathbf{x}^k\end{aligned}$$

- **ADMM iterates:**

$$\begin{aligned}\mathbf{x}^{k+1} &= (\mathbf{A}^\top \mathbf{A} + \rho \mathbf{I})^{-1} (\mathbf{A}^\top \mathbf{b} + \rho (\mathbf{z}^k - \mathbf{u}^k)) \\ \mathbf{z}^{k+1} &= \mathcal{S}_{\lambda/\rho} (\mathbf{x}^{k+1} + \mathbf{u}^k) & k = 0, 1, 2, \dots \\ \mathbf{u}^{k+1} &= \mathbf{u}^k + (\mathbf{x}^{k+1} - \mathbf{z}^{k+1})\end{aligned}$$

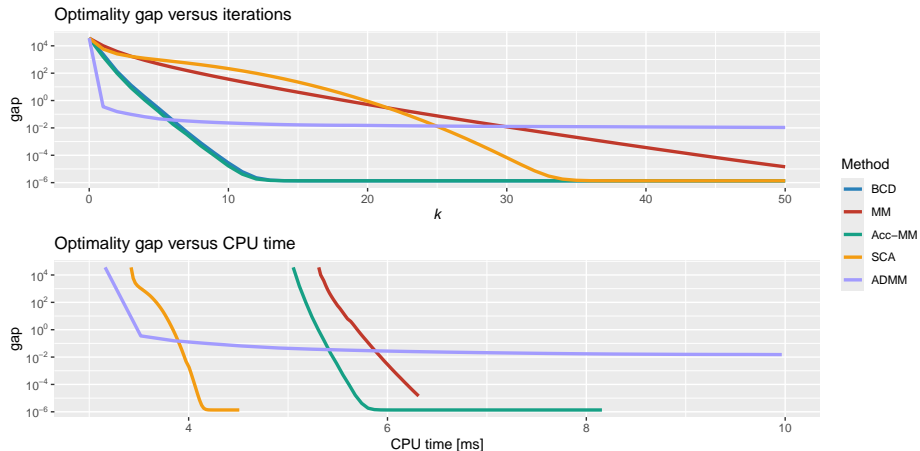
## Comparison of Methods

- The BCD method updates each element sequentially ( $n = 100$ ), which leads to a high computational cost (CPU time).
- The Jacobi method is the parallel version of BCD, is not guaranteed to converge, and is similar to SCA but lacks  $\tau$  and a smoothing step.
- The MM method requires computing the largest eigenvalue of  $\mathbf{A}^T \mathbf{A}$  and uses a conservative upper-bound  $\kappa$  for all elements.
- The SCA method uses  $\text{diag}(\mathbf{A}^T \mathbf{A})$  instead of a common  $\kappa$ , which leads to faster convergence due to element-specific updates.
- The ADMM method converges with lower accuracy but is often sufficient for practical applications.



# Numerical Comparison

Comparison of different iterative methods for the  $\ell_2 - \ell_1$ -norm minimization:



# Outline

- 1 Solvers
- 2 Gradient Methods
- 3 Interior-Point Methods (IPM)
- 4 Fractional Programming (FP) Methods
- 5 BCD
- 6 MM
- 7 SCA
- 8 ADMM
- 9 Numerical Comparison
- 10 Summary**

# Summary

- Solvers for convex and nonconvex problems are available in all programming languages, often used via modeling frameworks.
- Solvers use methods like gradient descent, Newton's method, and interior-point methods, but users typically don't need to understand these details.
- Advanced users may develop custom algorithms for specific problems, requiring more effort and knowledge, such as the Dinkelbach method or Charnes-Cooper-Schaible transform for fractional problems.
- Iterative algorithmic frameworks break complex problems into easier ones:
  - Bisection
  - Block Coordinate Descent (BCD)
  - Majorization-Minimization (MM)
  - Successive Convex Approximation (SCA)
  - Alternating Direction Method of Multipliers (ADMM)

- Beck, A. 2017. *First-Order Methods in Optimization*. MOS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics (SIAM).
- Bertsekas, D. P. 1999. *Nonlinear Programming*. Athena Scientific.
- Bertsekas, D. P., and J. N. Tsitsiklis. 1997. *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific.
- Boyd, S. P., and L. Vandenberghe. 2004. *Convex Optimization*. Cambridge University Press.
- Boyd, S., N. Parikh, E. Chu, B. Peleato, and J. Eckstein. 2010. *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*. Foundations and Trends in Machine Learning, Now Publishers.
- Fu, A., B. Narasimhan, and S. Boyd. 2020. “CVXR: An R Package for Disciplined Convex Optimization.” *Journal of Statistical Software* 94 (14): 1–34.

## References II

- Grant, M., and S. Boyd. 2008. “Graph Implementations for Nonsmooth Convex Programs.” In *Recent Advances in Learning and Control*, edited by V. Blondel, S. Boyd, and H. Kimura, 95–110. Lecture Notes in Control and Information Sciences. Springer-Verlag.
- . 2014. *CVX: Matlab Software for Disciplined Convex Programming*. <http://cvxr.com/cvx>.
- Grippo, L., and M. Sciandrone. 2000. “On the Convergence of the Block Nonlinear Gauss–Seidel Method Under Convex Constraints.” *Operations Research Letters* 26 (3): 127–36.
- Hunter, D. R., and K. Lange. 2004. “A Tutorial on MM Algorithms.” *The American Statistician* 58: 30–37.
- Kumar, S., J. Ying, J. V. M. Cardoso, and D. P. Palomar. 2019. “Structured Graph Learning via Laplacian Spectral Constraints.” In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*. Vancouver, Canada.

## References III

- . 2020. “A Unified Framework for Structured Graph Learning via Spectral Constraints.” *Journal of Machine Learning Research (JMLR)*, 1–60.
- Löfberg, J. 2004. “YALMIP: A Toolbox for Modeling and Optimization in MATLAB.” In *Proceedings of the CACSD Conference*. Taipei, Taiwan.
- Nemirovski, A. 2001. “Lectures on Modern Convex Optimization.” In *Society for Industrial and Applied Mathematics (SIAM)*.
- Nesterov, Y. 2018. *Lectures on Convex Optimization*. 2nd ed. Springer.
- Nesterov, Y., and A. Nemirovskii. 1994. *Interior-Point Polynomial Algorithms in Convex Programming*. Philadelphia, PA: SIAM.
- Nocedal, J., and S. J. Wright. 2006. *Numerical Optimization*. Springer Verlag.

- Palomar, D. P. 2025. *Portfolio Optimization: Theory and Application*. Cambridge University Press.
- Razaviyayn, M., M. Hong, and Z. Luo. 2013. “A Unified Convergence Analysis of Block Successive Minimization Methods for Nonsmooth Optimization.” *SIAM Journal on Optimization* 23 (2): 1126–53.
- Scutari, G., F. Facchinei, P. Song, D. P. Palomar, and J.-S. Pang. 2014. “Decomposition by Partial Linearization: Parallel Optimization of Multi-Agent Systems.” *IEEE Transactions on Signal Processing* 62 (3): 641–56.
- Scutari, G., and Y. Sun. 2018. “Parallel and Distributed Successive Convex Approximation Methods for Big-Data Optimization.” In *Multi-Agent Optimization*, edited by F. Facchinei and J. S. Pang, 141–308. Lecture Notes in Mathematics, Springer.

- Song, J., P. Babu, and D. P. Palomar. 2015. "Sparse Generalized Eigenvalue Problem via Smooth Optimization." *IEEE Transactions on Signal Processing* 63 (7): 1627–42.
- Sun, Y., P. Babu, and D. P. Palomar. 2017. "Majorization-Minimization Algorithms in Signal Processing, Communications, and Machine Learning." *IEEE Transactions on Signal Processing* 65 (3): 794–816.
- Varadhan, R., and C. Roland. 2008. "Simple and Globally Convergent Methods for Accelerating the Convergence of Any EM Algorithm." *Scandinavian Journal of Statistics* 35 (2): 335–53.
- Zibulevsky, M., and M. Elad. 2010. "L1 - L2 Optimization in Signal and Image Processing." *IEEE Signal Processing Magazine*, May, 76–88.