# Portfolio Optimization

## Optimization Algorithms

Daniel P. Palomar (2024). *Portfolio Optimization: Theory and Application.*
Cambridge University Press.

portfoliooptimizationbook.com

# Outline

# Abstract

Over the past century, the development of efficient algorithms for solving convex optimization problems has seen significant advancements. In 1947, Dantzig introduced the simplex method for linear programming (LP), which, despite its exponential worst-case complexity, became widely used. In 1984, Karmarkar's interior-point method revolutionized LP by offering polynomial time complexity. This innovation spurred further research, extending interior-point methods to quadratic programming (QP) and linear complementarity problems. In 1994, Nesterov and Nemirovskii advanced the field with the theory of self-concordant functions, enabling the application of log-barrier function-based algorithms to a broader range of convex problems, including semidefinite programming (SDP) and second-order cone programming (SOCP). Additionally, various specialized techniques like block-coordinate descent, majorization-minimization, and successive convex approximation have been developed to create customized algorithms for specific problems, often enhancing complexity and convergence rates. These slides will delve into a wide array of such practical algorithms (Palomar 2024, Appendix B).

# Outline

## Solvers

- 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.

## Popular solvers

- **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.

# Popular solvers

- **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 $O(n^3 L)$.
  - $n$: number of variables.
  - $L$: number of accuracy digits of the solution.

- **Specific complexities:**
  - **LP:** $O((m + n)^{3/2} n^2 L)$.
  - **QCQP:** $O(\sqrt{m}(m + n)n^2 L)$.
  - **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 $O(n^6 L)$.
  - Complexity for solving SOCP is higher than LP, QP, and QCQP; even higher for SDP.

## Interface with solvers

- **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:**

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

  - Solution depends on the choice of the norm.

- **Norm approximation with Euclidean or $\ell_2$-norm:**

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

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

## Interface with solvers

- **Norm approximation with Chebyshev or $\ell_\infty$-norm:**

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

- Rewritten as LP:

$$\underset{\boldsymbol{x},t}{\text{minimize}} \quad t$$
$$\text{subject to} \quad -t\boldsymbol{1} \leq \boldsymbol{Ax} - \boldsymbol{b} \leq t\boldsymbol{1}$$

- Equivalent form:

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

- Matlab code:

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

## Interface with solvers

- **Norm approximation problem with Manhattan or $\ell_1$-norm:**

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_1$$

- Rewritten as LP:

$$\underset{\boldsymbol{x},\boldsymbol{t}}{\text{minimize}} \quad \mathbf{1}^{\mathsf{T}}\boldsymbol{t}$$
$$\text{subject to} \quad -\boldsymbol{t} \leq \boldsymbol{A}\boldsymbol{x} - \boldsymbol{b} \leq \boldsymbol{t}$$

- Equivalent form:

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

- Matlab code:

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

## Interface with solvers

- **Euclidean norm approximation problem with linear constraints:**

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

- Equivalent form:

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

## Interface with solvers

- **Euclidean norm approximation problem with linear constraints: (cont'd)**
  - Equivalent form:

$$
\begin{aligned}
\underset{\boldsymbol{x},\boldsymbol{y},t,\boldsymbol{s}_l,\boldsymbol{s}_u}{\text{minimize}} \quad & t \\
\text{subject to} \quad & \boldsymbol{A}\boldsymbol{x} - \boldsymbol{b} = \boldsymbol{y} \\
& \boldsymbol{C}\boldsymbol{x} = \boldsymbol{d} \\
& \boldsymbol{x} - \boldsymbol{s}_l = \boldsymbol{l} \\
& \boldsymbol{x} + \boldsymbol{s}_u = \boldsymbol{u} \\
& \boldsymbol{s}_l, \boldsymbol{s}_u \geq \boldsymbol{0} \\
& \|\boldsymbol{y}\|_2 \leq t
\end{aligned}
$$

  - Equivalent form in matrix notation:

$$
\begin{aligned}
\underset{\boldsymbol{x},\boldsymbol{y},t,\boldsymbol{s}_l,\boldsymbol{s}_u}{\text{minimize}} \quad & \begin{bmatrix} \boldsymbol{0}^{\mathsf{T}} & \boldsymbol{0}^{\mathsf{T}} & \boldsymbol{0}^{\mathsf{T}} & \boldsymbol{0}^{\mathsf{T}} & 1 \end{bmatrix} \bar{\boldsymbol{x}} \\
\text{subject to} \quad & \begin{bmatrix} \boldsymbol{A} & & & & -\boldsymbol{I} \\ \boldsymbol{C} & & & & \\ \boldsymbol{I} & -\boldsymbol{I} & & & \\ \boldsymbol{I} & & & \boldsymbol{I} & \end{bmatrix} \bar{\boldsymbol{x}} \leq \begin{bmatrix} \boldsymbol{b} \\ \boldsymbol{d} \\ \boldsymbol{l} \\ \boldsymbol{u} \end{bmatrix} \\
& \bar{\boldsymbol{x}} \in \boldsymbol{R}^n \times \boldsymbol{R}_+^n \times \bar{\boldsymbol{R}}_+^n \times \boldsymbol{Q}^m
\end{aligned}
$$

# Interface with solvers

- **Euclidean norm approximation problem with linear constraints: (cont'd)**
  - **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 frameworks

- **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.
  - Supports integer constraints.
  - 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.

## Modeling frameworks

- **Example: Constrained Euclidean norm approximation in CVX:**
  - **Problem statement:**

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

  - **Matlab code:**

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

## Modeling frameworks

- **Example: Constrained Euclidean norm approximation in CVX: (cont'd)**
  - **R code:**

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

  - **Python code:**

```python
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

# Gradient methods

- **Unconstrained optimization problem:**

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

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

- **Iterative methods:**
  - Produce a sequence of iterates $x^0, x^1, x^2, \ldots$
  - Sequence may or may not converge to an optimal solution $x^\star$.

- **Ideal case with convex $f$:**
  - As iterations proceed ($k \to \infty$):
    - Objective function converges to the optimal value:

$$f\left(x^k\right) \to p^\star$$

    - Gradient tends to zero:

$$\nabla f\left(x^k\right) \to \mathbf{0}$$

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

## Descent methods

- **Descent methods (gradient methods):**
    - Satisfy the property: $f\left(\mathbf{x}^{k+1}\right) < f\left(\mathbf{x}^{k}\right)$.
    - 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:

    $$\nabla f\left(\mathbf{x}\right)^{\mathsf{T}}\mathbf{d} < 0$$

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

# Line search

- **Line search:**
  - Procedure to choose the stepsize $\alpha$.
  - Two widely used methods due to good theoretical convergence and practical performance:

- **Exact line search:**
  - Solves the univariate optimization problem:

$$\alpha = \underset{\alpha > 0}{\arg\min} \, f(\boldsymbol{x} + \alpha \boldsymbol{d}).$$

- **Backtracking line search (Armijo rule):**
  - Starting at $\alpha = 1$, repeat $\alpha \leftarrow \beta\alpha$ until:

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

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:

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

  which is a descent direction since $\nabla f\left(\boldsymbol{x}\right)^{\mathsf{T}} \boldsymbol{d} < 0$.

- **Gradient descent update:**

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

- **Stopping criterion:**
  - Common heuristic: $\|\nabla f(\boldsymbol{x})\|_2 \leq \epsilon$.

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

# Gradient descent method

## Gradient descent method

**Initialization:**

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

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

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

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

4. $k \leftarrow k + 1$

**Until:** convergence

# Newton's method

- **Newton's method:**
  - A descent method using both the gradient and the Hessian of $f$.
  - **Search direction:**
    $$\boldsymbol{d} = -\nabla^2 f(\boldsymbol{x})^{-1} \nabla f(\boldsymbol{x}),$$
    - Assumes $f$ is convex, twice continuously differentiable, and the Hessian matrix is positive definite for all $\boldsymbol{x}$.

- **Second-order approximation:**
  - $\boldsymbol{x} + \boldsymbol{d}$ minimizes the second-order approximation of $f(\boldsymbol{x})$ around $\boldsymbol{x}$:
    $$\hat{f}(\boldsymbol{x} + \boldsymbol{v}) = f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^\mathsf{T} \boldsymbol{v} + \frac{1}{2} \boldsymbol{v}^\mathsf{T} \nabla^2 f(\boldsymbol{x}) \boldsymbol{v}.$$

- **Newton's method update:**
  $$\boldsymbol{x}^{k+1} = \boldsymbol{x}^k - \alpha^k \nabla^2 f\left(\boldsymbol{x}^k\right)^{-1} \nabla f\left(\boldsymbol{x}^k\right).$$

# Newton's method

- **Newton decrement:**
  - Measures the proximity of **x** to an optimal point:

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

  - Estimates $f(\boldsymbol{x}) - p^\star$:

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

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

- **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

## Newton's method

**Initialization:**

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

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

1. Compute Newton direction and decrement:

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

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

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

4. $k \leftarrow k + 1$

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

# Convergence

- **Convergence of descent methods:**
  - Ideally, the sequence $\{\boldsymbol{x}^k\}$ 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 $\{\boldsymbol{x}^k\}$ 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 $\{\boldsymbol{x}^k\}$ 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: $\alpha^k \to 0$ 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:**

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f(\boldsymbol{x})$$
$$\text{subject to} \quad \boldsymbol{x} \in \mathcal{X},$$

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

- **Descent method:**
  - Iterative update:

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

  where $\boldsymbol{d}^k$ is a descent direction.
  - Potential issue: $\boldsymbol{x}^{k+1}$ 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):

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

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

# Projected gradient methods

- **Generalized gradient projection method:**
  - Iterative update:

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

    - $\mathbf{d}^k = \bar{\mathbf{x}}^k - \mathbf{x}^k$ 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:**

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

- $\bar{\mathbf{x}}^k$: Projection of $\mathbf{x}^k - s^k \nabla f \left( \mathbf{x}^k \right)$ 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 $\boldsymbol{x} \in \mathcal{X}$.
- **Constrained Newton's method:**
  - **Iterative update:**

$$\bar{\boldsymbol{x}}^k = \arg\min_{\boldsymbol{x} \in \mathcal{X}} \left\{ \nabla f\left(\boldsymbol{x}^k\right)^\mathsf{T} \left(\boldsymbol{x} - \boldsymbol{x}^k\right) + \frac{1}{2s^k} \left(\boldsymbol{x} - \boldsymbol{x}^k\right)^\mathsf{T} \nabla^2 f\left(\boldsymbol{x}^k\right) \left(\boldsymbol{x} - \boldsymbol{x}^k\right) \right\}$$

$$\boldsymbol{x}^{k+1} = \boldsymbol{x}^k + \alpha^k \left(\bar{\boldsymbol{x}}^k - \boldsymbol{x}^k\right).$$

  - $\bar{\boldsymbol{x}}^k$: 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 $\boldsymbol{x}^k$.
  - The main difficulty is solving the quadratic subproblem to find $\bar{\boldsymbol{x}}^k$.
    - 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

- **Convergence of gradient projection methods:**
  - Detailed in (Bertsekas 1999).
  - **Theorem: Convergence of gradient projection methods**
    - Suppose $\{\boldsymbol{x}^k\}$ 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 $\{\boldsymbol{x}^k\}$ 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

# 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:**

$$\begin{aligned} \underset{\boldsymbol{x}}{\text{minimize}} \quad & f_0(\boldsymbol{x}) \\ \text{subject to} \quad & f_i(\boldsymbol{x}) \leq 0, \qquad i = 1, \ldots, m \\ & \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b} \end{aligned}$$

- **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 $Ax = b$:**
  - Solutions can be expressed as:

$$\{x \in \mathbb{R}^n \mid Ax = b\} = \{Fz + x_0 \mid z \in \mathbb{R}^{n-p}\},$$

  - $x_0$ is any particular solution to $Ax = b$.
  - $F \in \mathbb{R}^{n \times (n-p)}$ spans the nullspace of $A$, i.e., $AF = 0$.
- **Reduced or eliminated problem:**
  - Equivalent to the original problem:

$$\begin{aligned}
\underset{z}{\text{minimize}} \quad & \tilde{f}_0(z) \triangleq f_0(Fz + x_0) \\
\text{subject to} \quad & \tilde{f}_i(z) \triangleq f_i(Fz + x_0) \leq 0, \qquad i = 1, \ldots, m,
\end{aligned}$$

  - Gradients and Hessians:

$$\nabla \tilde{f}_i(z) = F^T \nabla f_i(x)$$
$$\nabla^2 \tilde{f}_i(z) = F^T \nabla^2 f_i(x) F.$$

## Indicator function

- **Reformulation via indicator function:**
  - Equivalent form of the original problem:

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f_0(\boldsymbol{x}) + \sum_{i=1}^{m} I_-\left(f_i(\boldsymbol{x})\right)$$
$$\text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b},$$

  - Indicator function definition:

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

- **Characteristics of the reformulated problem:**
  - Inequality constraints are eliminated.
  - Indicator function is included in the objective.
  - Drawbacks:
    - Indicator function is noncontinuous.
    - Indicator function is nondifferentiable.
    - Not practical for optimization.

# Logarithmic barrier

- **Smooth approximation of the indicator function:**
    - Popular choice: **logarithmic barrier**:

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

    - Parameter $t > 0$ controls the approximation.
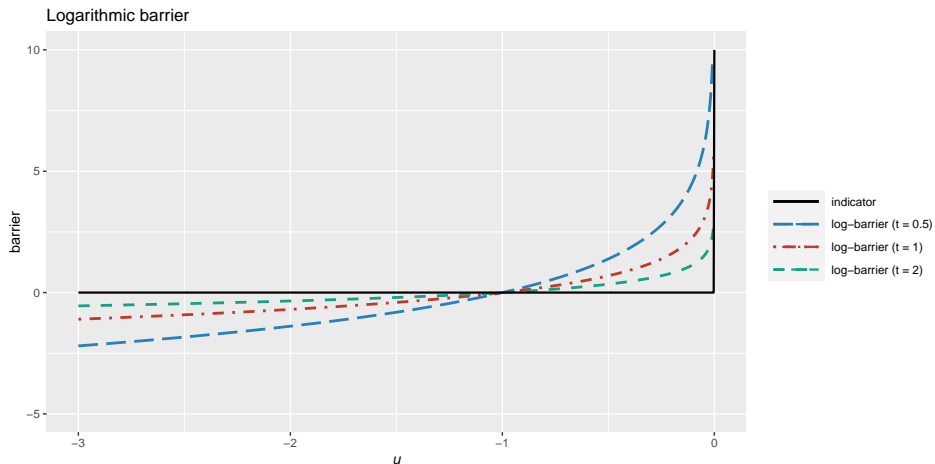    - Approximation improves as $t \to \infty$.

- **Approximate problem using the logarithmic barrier:**
    - Reformulated problem:

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

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



Logarithmic barrier

# Logarithmic barrier

- **Logarithmic barrier function:**
  - Overall barrier function (excluding the $1/t$ factor):

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

  which is convex (from composition rules).
  - Gradient and Hessian:

  $$\nabla\phi(\boldsymbol{x}) = \sum_{i=1}^{m} \frac{1}{-f_i(\boldsymbol{x})} \nabla f_i(\boldsymbol{x})$$

  $$\nabla^2\phi(\boldsymbol{x}) = \sum_{i=1}^{m} \frac{1}{f_i(\boldsymbol{x})^2} \nabla f_i(\boldsymbol{x}) \nabla f_i(\boldsymbol{x})^{\mathsf{T}} + \sum_{i=1}^{m} \frac{1}{-f_i(\boldsymbol{x})} \nabla^2 f_i(\boldsymbol{x}).$$

## Central path

- **Central path:** Defined as the curve $\{\boldsymbol{x}^\star(t) \mid t > 0\}$, where $\boldsymbol{x}^\star(t)$ is the solution to

$$\begin{aligned} \underset{\boldsymbol{x}}{\text{minimize}} \quad & tf_0(\boldsymbol{x}) + \phi(\boldsymbol{x}) \\ \text{subject to} \quad & \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b}, \end{aligned}$$

  which can be solved via Newton's method.

- **Solution to the central path problem:**
  - Ignoring equality constraints for simplicity:

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

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

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

# Central path

- **Convergence to optimal value:** $f_0(x^\star(t)) \to p^\star$ as $t \to \infty$.
  - From Lagrange duality theory:

$$p^\star \geq g\left(\lambda^\star(t)\right)$$
$$= L\left(x^\star(t); \lambda^\star(t)\right)$$
$$= f_0\left(x^\star(t)\right) - m/t.$$

- **Connection with KKT conditions:**
  - $x^\star(t)$ and $\lambda^\star(t)$ satisfy:

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

- Difference with original KKT conditions:
  - Complementary slackness is approximately satisfied.
  - Approximation improves as $t \to \infty$.

# Central path

Example of central path of an LP:

# Barrier method

- **Smooth approximation with logarithmic barrier:**
  - Log-barrier problem is a smooth approximation of the original problem.
  - Approximation improves as $t \to \infty$.
- **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.
    - 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 $x^\star(t)$ using Newton's method.
  - **Interior-point methods (IPM):**
    - Achieve this trade-off.
    - For each $t > 0$, $x^\star(t)$ is strictly feasible and lies in the interior of the feasible set.

# Barrier method

- **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 $f_0(\boldsymbol{x}) - p^\star \leq \epsilon$.
  - Refer to (S. P. Boyd and Vandenberghe 2004) for practical details.

# Barrier method

## Barrier method for constrained optimization

**Initialization:**

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

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

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

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

## Barrier method

- **Example: Barrier method for LP:**
  - Consider the LP:

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

  - Use the barrier method with different $\mu$ values.
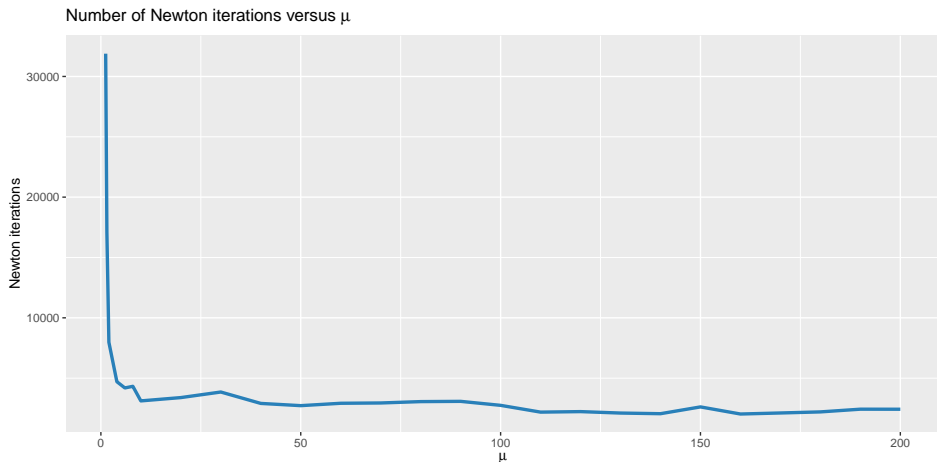  - **Convergence analysis:**
    - Case: $m = 100$ inequalities, $n = 50$ variables.
    - $\epsilon = 10^{-6}$ for the duality gap.
    - Centering problem solved via Newton's method.
  - **Observation:**
    - Total number of Newton iterations is not very sensitive to $\mu$ as long as $\mu \geq 10$.

# Barrier method

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



Number of Newton iterations versus $\mu$

# Convergence

- **Termination criterion:**
  - Number of outer iterations (centering steps) required:

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

  - Solving for $k$:

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

  - $\lceil \cdot \rceil$ 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:**
  - Requires a strictly feasible initial point $x^0$ (such that $f_i(x^0) < 0$).
  - 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:

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

  - 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:

$$\begin{aligned}
\underset{\bm{x}, s}{\text{minimize}} \quad & s \\
\text{subject to} \quad & f_i(\bm{x}) \le s, \qquad i = 1, \dots, m \\
& \bm{A}\bm{x} = \bm{b}.
\end{aligned}$$

  - **Constructing a strictly feasible point:**
    - Choose any $\bm{x}$ that satisfies the equality constraints.
    - Choose $s$ such that $s > f_i(\bm{x})$, e.g., $s = 1.1 \times \max_i \{f_i(\bm{x})\}$.
    - This provides an initial strictly feasible point for the Phase I problem.

- **Solving the Phase I problem:**
  - Obtain $(\bm{x}^\star, s^\star)$.
  - If $s^\star < 0$:
    - $\bm{x}^\star$ is a strictly feasible point.
    - Can be used in the barrier method to solve the original problem.
  - If $s^\star > 0$:
    - No feasible point exists.
    - No need to attempt solving 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

# Fractional programming (FP) methods

- **Concave-convex fractional program (FP):**

$$\underset{x}{\text{maximize}} \quad \frac{f(x)}{g(x)}$$
$$\text{subject to} \quad x \in \mathcal{X},$$

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

- **Nature of FPs:**
  - Nonconvex problems, generally difficult to solve
  - Concave-convex FP is a quasiconvex optimization problem, making it more tractable

- **Methods to solve concave-convex FP:**
  - **Iterative bisection method**
  - **Dinkelbach method**
  - **Schaible transform**

# FP: Bisection method

- **Problem reformulation:**
  - Solve a sequence of convex feasibility problems:

$$
\begin{aligned}
\underset{\boldsymbol{x}}{\text{find}} \quad & \boldsymbol{x} \\
\text{subject to} \quad & tg(\boldsymbol{x}) \leq f(\boldsymbol{x}) \\
& \boldsymbol{x} \in \mathcal{X}
\end{aligned}
$$

  - $t > 0$ 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

# FP: Bisection method

- Starts with an interval $[l, u]$ known to contain the optimal value $p^\star$ 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 $p^\star \in [l, u]$.

**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 $l = t$.

# FP: Dinkelbach method

- **Dinkelbach transform:**
    - **Objective:** Reformulate the original concave-convex FP into a sequence of simpler convex problems
    - **Reformulated problem:**

$$\underset{\boldsymbol{x}}{\text{maximize}} \quad f(\boldsymbol{x}) - y^k g(\boldsymbol{x})$$
$$\text{subject to} \quad \boldsymbol{x} \in \mathcal{X}$$

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

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

# FP: Dinkelbach method

- 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 $x^0$.
- Set $k \leftarrow 0$.

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

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

**Until:** convergence

# FP: Charnes-Cooper transform

- **Linear fractional program (LFP):**

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \dfrac{c^\mathsf{T} x + d}{e^\mathsf{T} x + f} \\ \text{subject to} & Gx \leq h \\ & Ax = b \end{array}$$

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

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

$$\begin{array}{ll} \underset{y,t}{\text{minimize}} & c^\mathsf{T} y + dt \\ \text{subject to} & Gy \leq ht \\ & Ay = bt \\ & e^\mathsf{T} y + ft = 1 \\ & t \geq 0 \end{array}$$

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

# FP: Charnes-Cooper transform

**Proof:**

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

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

# FP: Schaible transform

- **Concave-convex fractional program (FP):**

$$\underset{x}{\text{maximize}} \quad \frac{f(\pmb{x})}{g(\pmb{x})}$$
$$\text{subject to} \quad \pmb{x} \in \mathcal{X}$$

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

$$\underset{y,t}{\text{maximize}} \quad tf\left(\frac{\pmb{y}}{t}\right)$$
$$\text{subject to} \quad tg\left(\frac{\pmb{y}}{t}\right) \leq 1$$
$$t \geq 0$$
$$\pmb{y}/t \in \mathcal{X}$$

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

# FP: Schaible transform

**Proof:**

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

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

# Outline

# BCD

- **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:**

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

where $f$ is the (possibly nonconvex) objective function and each $\mathcal{X}_i$ is a convex set.
  - **Partitioning:** Variable $\mathbf{x}$ is partitioned into $n$ blocks $\mathbf{x} = (\mathbf{x}_1, \ldots, \mathbf{x}_n)$
- **Method description:**
  - **Iterative process:** Produces a sequence of iterates $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \ldots$ that converge to $\mathbf{x}^\star$
  - **Update rule:** Optimize the problem with respect to each block $\mathbf{x}_i$ sequentially
  - **Inner iterations:** At each outer iteration $k$, execute $n$ inner iterations sequentially:

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

# BCD

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

## BCD for separable problems

**Initialization:**

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

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

1. Execute $n$ inner iterations sequentially:

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

2. $k \leftarrow k + 1$

**Until:** convergence

## BCD: Convergence

- BCD enjoys monotonicity, i.e., $f\left(\mathbf{x}^{k+1}\right) \leq f\left(\mathbf{x}^k\right)$
- **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):**
  - **Objective:** Execute $n$ inner iterations in parallel instead of sequentially
  - **Update rule:**

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

- **Jacobi method:**
  - **Description:** Parallel update of block variables
  - **Algorithmic attractiveness:** Potentially faster due to parallel execution

- **Convergence properties:**
  - **Issue:** Jacobi method does not enjoy nice convergence properties
  - **Condition for convergence:** 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:**

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

- **Solution:**

$$x = \frac{1}{\|\boldsymbol{a}\|_2^2}\text{sign}\left(\boldsymbol{a}^\mathsf{T}\boldsymbol{b}\right)\left(|\boldsymbol{a}^\mathsf{T}\boldsymbol{b}| - \lambda\right)^+$$

  - Sign function:

$$\text{sign}(u) = \left\{ \begin{array}{rl} +1 & u > 0 \\ 0 & u = 0 \\ -1 & u < 0 \end{array} \right.$$

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

# BCD example: Soft-thresholding operator

- **Compact form:**

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

- **Soft-thresholding operator:**

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

# BCD example: $\ell_2 - \ell_1$-norm minimization

- **Problem formulation:**

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

- **Solution approach:**
    - **Standard method:** Can be solved with a QP solver
    - **Iterative algorithm via BCD:** via soft-thresholding operator (Zibulevsky and Elad 2010)
- **BCD method:**
    - **Variable partitioning:** Divide the variable into each constituent element $\boldsymbol{x} = (x_1, \ldots, x_n)$
    - **Sequence of problems** at each iteration $k = 0, 1, 2, \ldots$ for each element $i = 1, \ldots, n$:

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

    where $\tilde{\boldsymbol{b}}_i^k \triangleq \boldsymbol{b} - \sum_{j<i} \boldsymbol{a}_j x_j^{k+1} - \sum_{j>i} \boldsymbol{a}_j x_j^k$.
- **Iterative algorithm:** For $k = 0, 1, 2, \ldots$:

$$x_i^{k+1} = \frac{1}{\|\boldsymbol{a}_i\|_2^2} \mathcal{S}_\lambda \left( \boldsymbol{a}_i^\mathsf{T} \tilde{\boldsymbol{b}}_i^k \right), \quad i = 1, \ldots, n$$

# BCD example: $\ell_2 - \ell_1$-norm minimization

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



Optimality gap versus iterations

# Outline

## MM

- **Majorization-minimization (MM) method:**
  - **Objective:** Approximate a difficult optimization problem by a sequence of simpler problems.
  - **References:**
    - Concise tutorial: (Hunter and Lange 2004)
    - Long tutorial with applications: (Sun, Babu, and Palomar 2017)
    - Convergence analysis: (Razaviyayn, Hong, and Luo 2013)

- **Original problem:**

$$\begin{aligned} \underset{\boldsymbol{x}}{\text{minimize}} \quad & f(\boldsymbol{x}) \\ \text{subject to} \quad & \boldsymbol{x} \in \mathcal{X} \end{aligned}$$

where

- $f$ is the (possibly nonconvex) objective function
- $\mathcal{X}$ is a (possibly nonconvex) set.

# MM

- **MM method:**
  - **Iterative process:** Produces a sequence of iterates $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \ldots$ that converge to $\mathbf{x}^\star$.
  - **Surrogate function:** At iteration $k$, approximate $f(\mathbf{x})$ by a surrogate function $u\left(\mathbf{x}; \mathbf{x}^k\right)$ around the current point $\mathbf{x}^k$.
  - **Sequence of problems:**

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

# MM

Illustration of sequence of surrogate problems in MM:

# MM: Convergence

- **Conditions for surrogate function $u\left(x; x^k\right)$:**
  - **Upper-bound property:** $u\left(x; x^k\right) \geq f\left(x\right)$
  - **Touching property:** $u\left(x^k; x^k\right) = f\left(x^k\right)$
  - **Tangent property:** $u\left(x; x^k\right)$ must be differentiable with $\nabla u\left(x; x^k\right) = \nabla f\left(x\right)$

- **Consequences:**
  - **Monotonicity:** $f\left(x^{k+1}\right) \leq f\left(x^k\right)$
  - **Convergence:** If $\mathcal{X}$ is convex, every limit point of the sequence $\{x^k\}$ is a stationary point of the original problem

- **Majorizer construction:**
  - **Objective:** Find an appropriate majorizer $u\left(x; x^k\right)$ that satisfies the technical conditions and leads to a simpler surrogate problem
  - **Techniques and examples:** Refer to (Sun, Babu, and Palomar 2017)

# 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 $u\left(\mathbf{x}; \mathbf{x}^k\right)$.
2. Obtain next iterate by solving the majorized problem:

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

3. $k \leftarrow k + 1$

**Until:** convergence

# MM: Convergence

- **Versatility of MM framework:**
  - **Objective:** Derive practical algorithms
  - **Theoretical guarantees:** Convergence properties are well-established
- **Assumptions:**
  - Majorizer $u\left(\boldsymbol{x}; \boldsymbol{x}^k\right)$ satisfies the technical conditions.
  - Feasible set $\mathcal{X}$ is convex.
- **Convergence:** Every limit point of the sequence $\{\boldsymbol{x}^k\}$ is a stationary point of the original problem
- **Nonconvex feasible set $\mathcal{X}$:**
  - Convergence must be studied on a case-by-case basis.
  - Examples: (Song, Babu, and Palomar 2015; Sun, Babu, and Palomar 2017; Kumar et al. 2019, 2020).

- **MM convergence speed:**
  - **Issue:** MM may require many iterations to converge if the surrogate function $u\left(\boldsymbol{x}; \boldsymbol{x}^k\right)$ 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:**

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

  where the parameters are
  - $b \in \mathbb{R}_+^m$ (nonnegative elements)
  - $A \in \mathbb{R}_{++}^{m \times n}$ (positive elements).

- **Conventional LS solution:**
  - Not applicable due to nonnegativity constraints: $x^\star = (A^\mathsf{T} A)^{-1} A^\mathsf{T} b$.

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

# MM example: Nonnegative LS

- **Objective function:** $f(\boldsymbol{x}) = \frac{1}{2}\|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2$
- **Majorizer:**

$$u\left(\boldsymbol{x}; \boldsymbol{x}^k\right) = f\left(\boldsymbol{x}^k\right) + \nabla f\left(\boldsymbol{x}^k\right)^{\mathsf{T}}\left(\boldsymbol{x} - \boldsymbol{x}^k\right) + \frac{1}{2}\left(\boldsymbol{x} - \boldsymbol{x}^k\right)^{\mathsf{T}}\boldsymbol{\Phi}\left(\boldsymbol{x}^k\right)\left(\boldsymbol{x} - \boldsymbol{x}^k\right)$$

  - **Gradient:** $\nabla f\left(\boldsymbol{x}^k\right) = \boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x}^k - \boldsymbol{A}^{\mathsf{T}}\boldsymbol{b}$
  - **Matrix $\boldsymbol{\Phi}$:** $\boldsymbol{\Phi}\left(\boldsymbol{x}^k\right) = \text{Diag}\left(\frac{\left[\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x}^k\right]_1}{x_1^k}, \ldots, \frac{\left[\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x}^k\right]_n}{x_n^k}\right)$

- **Verification of majorizer properties:**
  - **Upper-bound property:** $u\left(\boldsymbol{x}; \boldsymbol{x}^k\right) \geq f\left(\boldsymbol{x}\right)$ (proved using Jensen's inequality)
  - **Touching property:** $u\left(\boldsymbol{x}^k; \boldsymbol{x}^k\right) = f\left(\boldsymbol{x}^k\right)$
  - **Tangent property:** $\nabla u\left(\boldsymbol{x}^k; \boldsymbol{x}^k\right) = \nabla f\left(\boldsymbol{x}^k\right)$

## MM example: Nonnegative LS

- **Sequence of majorized problems:**

$$\underset{\boldsymbol{x} \geq \boldsymbol{0}}{\text{minimize}} \quad \nabla f\left(\boldsymbol{x}^k\right)^{\mathsf{T}} \boldsymbol{x} + \tfrac{1}{2} \left(\boldsymbol{x} - \boldsymbol{x}^k\right)^{\mathsf{T}} \boldsymbol{\Phi}\left(\boldsymbol{x}^k\right) \left(\boldsymbol{x} - \boldsymbol{x}^k\right)$$

- **Solution:** $\boldsymbol{x} = \boldsymbol{x}^k - \boldsymbol{\Phi}\left(\boldsymbol{x}^k\right)^{-1} \nabla f\left(\boldsymbol{x}^k\right)$

- **Iterative update:**

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

where $c_i^k = \frac{[\boldsymbol{A}^{\mathsf{T}}\boldsymbol{b}]_i}{[\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x}^k]_i}$ and $\odot$ denotes elementwise product.

# MM example: Nonnegative LS

Convergence of MM for the nonnegative LS:



Optimality gap vs iterations

# 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:**

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$$
$$\text{subject to} \quad \boldsymbol{x}_i \in \mathcal{X}_i, \quad i = 1, \ldots, n$$

  - **Partitioning:** Variables are partitioned into $n$ blocks $\boldsymbol{x} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$.
  - **Constraints:** Each block $\boldsymbol{x}_i$ is separately constrained.
- **Idea:** Solve the problem block by block as in BCD, but majorize each block $f(\boldsymbol{x}_i)$ with a surrogate function $u(\boldsymbol{x}_i; \boldsymbol{x}^k)$.
- **References:** (Razaviyayn, Hong, and Luo 2013) (Sun, Babu, and Palomar 2017).

# Block MM Procedure

1. **Initialization:** Start with an initial guess $x^0 = (x_1^0, \dots, x_n^0)$
2. **Iterative process:** For each outer iteration $k = 0, 1, 2, \dots$
   - **For each block** $i = 1, \dots, n$:
     - **Majorize:** Construct a surrogate function $u\left(x_i; x^k\right)$ for the block $f(x_i)$
     - **Update:** Solve the majorized problem for the block:

       $$x_i^{k+1} = \underset{x_i \in \mathcal{X}_i}{\arg\min}\, u\left(x_i; x^k\right)$$

   - **Update the full variable:** $x^{k+1} = (x_1^{k+1}, \dots, x_n^{k+1})$

# Outline

- **Successive convex approximation (SCA) method:**
  - Approximates a difficult optimization problem by a sequence of simpler convex problems.
  - Produces a sequence of iterates $x^0, x^1, x^2, \ldots$ that converge to $x^\star$.

- **Problem formulation:**

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

where

- $f$ is a (possibly nonconvex) objective function
- $\mathcal{X}$ is a convex set (nonconvex sets can be accommodated with more complexity).

# SCA

- **Iteration process:**
  - At iteration $k$, approximate $f(\mathbf{x})$ by a surrogate function $\tilde{f}\left(\mathbf{x}; \mathbf{x}^k\right)$ around $\mathbf{x}^k$.
  - Solve the sequence of simpler problems:

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

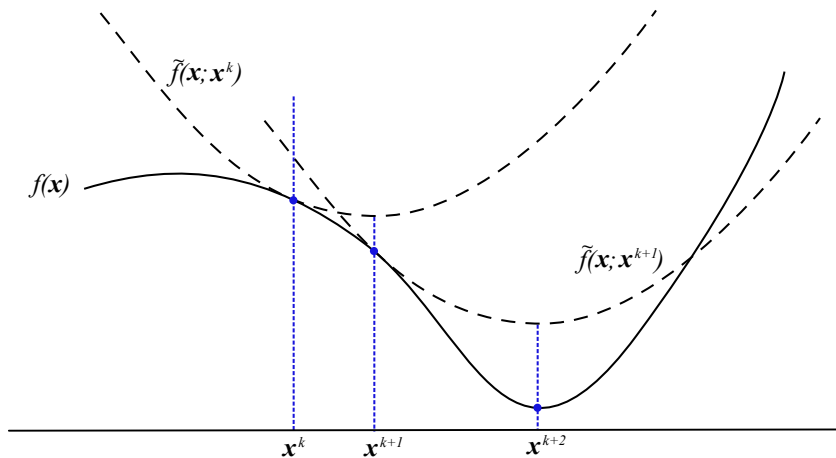  - Introduce a smoothing step to avoid oscillations:

  $$\hat{\mathbf{x}}^{k+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \tilde{f}\left(\mathbf{x}; \mathbf{x}^k\right)$$
  $$\mathbf{x}^{k+1} = \mathbf{x}^k + \gamma^k \left(\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k\right) \qquad k = 0, 1, 2, \ldots,$$

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

Illustration of sequence of surrogate problems in SCA:

# SCA: Convergence

- **Conditions for surrogate function $\tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right)$:**
  - Must be strongly convex on the feasible set $\mathcal{X}$.
  - Must be differentiable with $\nabla \tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right) = \nabla f\left(\boldsymbol{x}\right)$.
- **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*: Practical choice satisfying $\sum_{k=1}^{\infty} \gamma^k = +\infty$ and $\sum_{k=1}^{\infty} (\gamma^k)^2 < +\infty$.
    - Example 1: $\gamma^{k+1} = \gamma^k \left(1 - \epsilon\gamma^k\right)$, $\gamma^0 < 1/\epsilon$, $\epsilon \in (0, 1)$.
    - 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 $\alpha^k/\beta^k \to 0$.
- **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

## SCA algorithm

**Initialization:**

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

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

1. Construct surrogate of $f(\boldsymbol{x})$ around current point $\boldsymbol{x}^k$ as $\tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right)$.

2. Obtain intermediate point by solving the surrogate convex problem:

$$\hat{\boldsymbol{x}}^{k+1} = \underset{\boldsymbol{x} \in \mathcal{X}}{\arg \min} \, \tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right).$$

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

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

4. $k \leftarrow k + 1$

**Until:** convergence

# Gradient descent method as SCA

- **Unconstrained problem:**

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

- **SCA with surrogate function:**
  - Surrogate function:

$$\tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right) = f\left(\boldsymbol{x}^k\right) + \nabla f\left(\boldsymbol{x}^k\right)^{\mathsf{T}}\left(\boldsymbol{x} - \boldsymbol{x}^k\right) + \frac{1}{2\alpha^k}\|\boldsymbol{x} - \boldsymbol{x}^k\|^2$$

- **Minimizing the surrogate function:**
  - Set the gradient of $\tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right)$ to zero:

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

  - Solve for $\boldsymbol{x}$:

$$\boldsymbol{x} = \boldsymbol{x}^k - \alpha^k \nabla f\left(\boldsymbol{x}^k\right)$$

- **Iteration process:**
  - Update rule coincides with the gradient descent method:

$$\boldsymbol{x}^{k+1} = \boldsymbol{x}^k - \alpha^k \nabla f\left(\boldsymbol{x}^k\right), \qquad k = 0, 1, 2, \ldots$$

# Newton method as SCA

- **Including second-order information:**
  - Surrogate function with Hessian:

$$\tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right) = f\left(\boldsymbol{x}^k\right) + \nabla f\left(\boldsymbol{x}^k\right)^\mathsf{T}\left(\boldsymbol{x} - \boldsymbol{x}^k\right) + \frac{1}{2\alpha^k}\left(\boldsymbol{x} - \boldsymbol{x}^k\right)^\mathsf{T}\nabla^2 f\left(\boldsymbol{x}^k\right)\left(\boldsymbol{x} - \boldsymbol{x}^k\right)$$

- **Minimizing the surrogate function:**
  - Set the gradient of $\tilde{f}\left(\boldsymbol{x}; \boldsymbol{x}^k\right)$ to zero:

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

  - Solve for $\boldsymbol{x}$:

$$\boldsymbol{x} = \boldsymbol{x}^k - \alpha^k\nabla^2 f\left(\boldsymbol{x}^k\right)^{-1}\nabla f\left(\boldsymbol{x}^k\right)$$

- **Iteration process:**
  - Update rule coincides with Newton's method:

$$\boldsymbol{x}^{k+1} = \boldsymbol{x}^k - \alpha^k\nabla^2 f\left(\boldsymbol{x}^k\right)^{-1}\nabla f\left(\boldsymbol{x}^k\right), \qquad k = 0, 1, 2, \ldots$$

# Parallel SCA

- **Partitioned variables in SCA:**

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

  where variables are partitioned into $n$ separate blocks: $\boldsymbol{x} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$.

- **Parallel updates in SCA:**
  - Unlike BCD or MM, SCA updates variables in parallel with surrogate functions $\tilde{f}_i\left(\boldsymbol{x}_i; \boldsymbol{x}^k\right)$.
  - Update process for each block $i$:

$$\begin{aligned} \hat{\boldsymbol{x}}_i^{k+1} &= \underset{\boldsymbol{x}_i \in \mathcal{X}_i}{\arg \min} \; \tilde{f}_i\left(\boldsymbol{x}_i; \boldsymbol{x}^k\right) \\ \boldsymbol{x}_i^{k+1} &= \boldsymbol{x}_i^k + \gamma^k \left(\hat{\boldsymbol{x}}_i^{k+1} - \boldsymbol{x}_i^k\right) \end{aligned} \qquad i = 1, \ldots, n, \quad k = 0, 1, 2, \ldots$$

  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.

# SCA: Convergence

- **Technical conditions for surrogate function:**
  - Must be strongly convex on the feasible set $\mathcal{X}$.
  - Must be differentiable with $\nabla \tilde{f}\left(\mathbf{x}; \mathbf{x}^k\right) = \nabla f\left(\mathbf{x}\right)$.

- **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*: Practical choice 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 $\tilde{f}\left(\mathbf{x}; \mathbf{x}^k\right)$ (or each $\tilde{f}_i\left(\mathbf{x}_i; \mathbf{x}^k\right)$ 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:**

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

- **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:**
  - Partition variable $\boldsymbol{x}$ into elements $(x_1, \ldots, x_n)$.
  - Surrogate functions:

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

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

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

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

# SCA example: $\ell_2 - \ell_1$-norm minimization

- **SCA iterative algorithm:**
  - Update rule:

$$\hat{x}_i^{k+1} = \frac{1}{\tau + \|\boldsymbol{a}_i\|^2} \mathcal{S}_\lambda \left( \boldsymbol{a}_i^\mathsf{T} \tilde{\boldsymbol{b}}_i^k + \tau x_i^k \right)$$
$$x_i^{k+1} = x_i^k + \gamma^k \left( \hat{x}_i^{k+1} - x_i^k \right) \qquad i = 1, \ldots, n, \quad k = 0, 1, 2, \ldots$$

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

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

# SCA example: $\ell_2 - \ell_1$-norm minimization

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



Optimality gap vs iterations

# SCA example: Dictionary learning

- **Dictionary learning problem:**

$$\underset{\boldsymbol{D}, \boldsymbol{X}}{\text{minimize}} \quad \frac{1}{2}\|\boldsymbol{Y} - \boldsymbol{DX}\|_F^2 + \lambda\|\boldsymbol{X}\|_1$$
$$\text{subject to} \quad \|[\boldsymbol{D}]_{:,i}\| \leq 1, \qquad i = 1, \ldots, m.$$

  - $\|\boldsymbol{D}\|_F$: Frobenius norm of $\boldsymbol{D}$
  - $\|\boldsymbol{X}\|_1$: elementwise $\ell_1$-norm of $\boldsymbol{X}$

- **Matrix definitions:**
  - $\boldsymbol{D}$: dictionary matrix (fat matrix with columns explaining the columns of $\boldsymbol{Y}$)
  - $\boldsymbol{X}$: sparse matrix selecting a few columns of the dictionary

- **Bi-convex nature:**
  - Problem is not jointly convex in $(\boldsymbol{D}, \boldsymbol{X})$, but it is bi-convex.
  - For fixed $\boldsymbol{D}$, the problem is convex in $\boldsymbol{X}$.
  - For fixed $\boldsymbol{X}$, the problem is convex in $\boldsymbol{D}$.

- **Solution methods:**
  - BCD: updates $\boldsymbol{D}$ and $\boldsymbol{X}$ sequentially.
  - SCA: allows parallel updates of $\boldsymbol{D}$ and $\boldsymbol{X}$.

## SCA example: Dictionary learning

- **SCA approach:**
  - Surrogate functions:

$$\tilde{f}_1\left(\boldsymbol{D}; \boldsymbol{X}^k\right) = \frac{1}{2}\|\boldsymbol{Y} - \boldsymbol{D}\boldsymbol{X}^k\|_F^2$$

$$\tilde{f}_2\left(\boldsymbol{X}; \boldsymbol{D}^k\right) = \frac{1}{2}\|\boldsymbol{Y} - \boldsymbol{D}^k\boldsymbol{X}\|_F^2$$

- **Resulting convex problems:**
  - **Normalized least squares (LS) problem:**

$$\begin{array}{ll} \underset{\boldsymbol{D}}{\text{minimize}} & \frac{1}{2}\|\boldsymbol{Y} - \boldsymbol{D}\boldsymbol{X}^k\|_F^2 \\ \text{subject to} & \|[\boldsymbol{D}]_{:,i}\| \leq 1, \qquad i = 1, \ldots, m \end{array}$$

  - **Matrix version of the $\ell_2 - \ell_1$-norm problem:**

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

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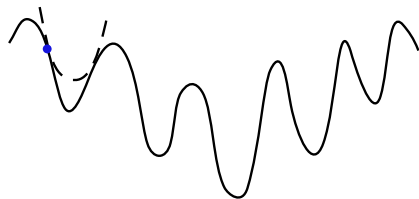
# MM versus SCA

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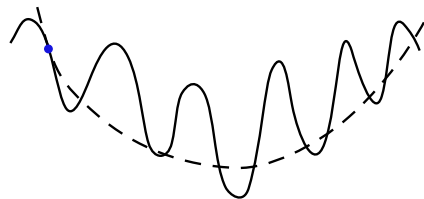  - **MM (Majorization-Minimization):**
    - Requires the surrogate function to be a global upper bound.
    - The surrogate function need not be convex.
    - Can be difficult to derive and too restrictive in some cases.
  - **SCA (Successive Convex Approximation):**
    - Relaxes the upper-bound condition.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



MM                                                          SCA

# MM versus 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 $\mathcal{X}$.
        - 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 $\boldsymbol{x} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$.
    - **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

# ADMM

- **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:**

$$\underset{x,z}{\text{minimize}} \quad f(x) + g(z)$$
$$\text{subject to} \quad Ax + Bz = c,$$

- Observe that the variables $x$ and $z$ are coupled via the constraint $Ax + Bz = c$.

## First Attempt: Dual ascent method

First attempt to decouple the variables.

**Dual ascent method:**

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

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

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

$$\boldsymbol{x}^{k+1} = \arg \min_{\boldsymbol{x}} \ f(\boldsymbol{x}) + (\boldsymbol{y}^k)^\mathsf{T} \boldsymbol{A}\boldsymbol{x}$$
$$\boldsymbol{z}^{k+1} = \arg \min_{\boldsymbol{z}} \ g(\boldsymbol{z}) + (\boldsymbol{y}^k)^\mathsf{T} \boldsymbol{B}\boldsymbol{z} \qquad k = 0, 1, 2, \ldots$$
$$\boldsymbol{y}^{k+1} = \boldsymbol{y}^k + \alpha^k \left( \boldsymbol{A}\boldsymbol{x}^{k+1} + \boldsymbol{B}\boldsymbol{z}^{k+1} - \boldsymbol{c} \right)$$

- Requires many technical assumptions and is often slow.

## Second Attempt: Method of multipliers

Second attempt to decouple the variables.

**Method of multipliers:**

- Uses augmented Lagrangian:

$$L_\rho(\boldsymbol{x}, \boldsymbol{z}; \boldsymbol{y}) \triangleq f(\boldsymbol{x}) + g(\boldsymbol{z}) + \boldsymbol{y}^\mathsf{T}(\boldsymbol{A}\boldsymbol{x} + \boldsymbol{B}\boldsymbol{z} - \boldsymbol{c}) + \frac{\rho}{2}\|\boldsymbol{A}\boldsymbol{x} + \boldsymbol{B}\boldsymbol{z} - \boldsymbol{c}\|_2^2$$

- Algorithm:

$$\begin{aligned}
\left(\boldsymbol{x}^{k+1}, \boldsymbol{z}^{k+1}\right) &= \arg\min_{\boldsymbol{x},\boldsymbol{z}} \ L_\rho(\boldsymbol{x}, \boldsymbol{z}; \boldsymbol{y}^k) \\
\boldsymbol{y}^{k+1} &= \boldsymbol{y}^k + \rho\left(\boldsymbol{A}\boldsymbol{x}^{k+1} + \boldsymbol{B}\boldsymbol{z}^{k+1} - \boldsymbol{c}\right)
\end{aligned} \qquad k = 0, 1, 2, \ldots$$

- Converges under more relaxed conditions but cannot decouple $\boldsymbol{x}$ and $\boldsymbol{z}$ due to $\|\boldsymbol{A}\boldsymbol{x} + \boldsymbol{B}\boldsymbol{z} - \boldsymbol{c}\|_2^2$ term.

## Third Attempt: ADMM

Third an final attempt to decouple the variables.

**ADMM:**

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

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

- Successfully decouples primal variables $\mathbf{x}$ and $\mathbf{z}$.
- Faster convergence with fewer technical conditions.
- Common to express ADMM updates using scaled dual variable $\mathbf{u}^k = \mathbf{y}^k / \rho$ as in the next algorithm.

# ADMM

## ADMM algorithm

**Initialization:**

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

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

1. Iterate primal and dual variables:

$$\boldsymbol{x}^{k+1} = \arg\min_{\boldsymbol{x}} \ f(\boldsymbol{x}) + \frac{\rho}{2} \left\| \boldsymbol{A}\boldsymbol{x} + \boldsymbol{B}\boldsymbol{z}^k - \boldsymbol{c} + \boldsymbol{u}^k \right\|_2^2$$

$$\boldsymbol{z}^{k+1} = \arg\min_{\boldsymbol{z}} \ g(\boldsymbol{z}) + \frac{\rho}{2} \left\| \boldsymbol{A}\boldsymbol{x}^{k+1} + \boldsymbol{B}\boldsymbol{z} - \boldsymbol{c} + \boldsymbol{u}^k \right\|_2^2$$

$$\boldsymbol{u}^{k+1} = \boldsymbol{u}^k + \left( \boldsymbol{A}\boldsymbol{x}^{k+1} + \boldsymbol{B}\boldsymbol{z}^{k+1} - \boldsymbol{c} \right);$$

2. $k \leftarrow k + 1$

**Until:** convergence

## ADMM: Convergence

- **Assumptions:**
  - $f(\boldsymbol{x})$ and $g(\boldsymbol{z})$ are convex.
  - Both the $\boldsymbol{x}$-update and the $\boldsymbol{z}$-update are solvable.
  - The Lagrangian has a saddle point.
- **Convergence of ADMM:**
  - **Residual convergence:** $\boldsymbol{Ax}^k + \boldsymbol{Bz}^k - \boldsymbol{c} \rightarrow \boldsymbol{0}$ as $k \rightarrow \infty$
    - Iterates approach feasibility.
  - **Objective convergence:** $f(\boldsymbol{x}) + g(\boldsymbol{z}) \rightarrow p^\star$ as $k \rightarrow \infty$
    - Objective function of the iterates approaches the optimal value.
  - **Dual variable convergence:** $\boldsymbol{y}^k \rightarrow \boldsymbol{y}^\star$ as $k \rightarrow \infty$
  - Detailed analysis in (S. Boyd et al. 2010) and references therein.
- **Practical considerations:**
  - $\{\boldsymbol{x}^k\}$ and $\{\boldsymbol{z}^k\}$ need not converge to optimal values without additional assumptions.
  - ADMM can be slow to converge to high accuracy.
  - Often converges to modest accuracy within a few tens of iterations, which is sufficient for many practical applications.
  - Different from the fast convergence of Newton's method.

# ADMM example: Constrained convex optimization

- **Generic convex optimization problem:**

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

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 $\mathcal{X}$:

$$g(\boldsymbol{x}) \triangleq \left\{ \begin{array}{ll} 0 & \boldsymbol{x} \in \mathcal{X} \\ +\infty & \text{otherwise}, \end{array} \right.$$

  - Formulate the equivalent problem:

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

## ADMM example: Constrained convex optimization

- **ADMM algorithm for the transformed problem:**
  - Update rules:

$$\boldsymbol{x}^{k+1} = \arg \min_{\boldsymbol{x}} \; f(\boldsymbol{x}) + \frac{\rho}{2} \left\| \boldsymbol{x} - \boldsymbol{z}^k + \boldsymbol{u}^k \right\|_2^2$$
$$\boldsymbol{z}^{k+1} = \left[ \boldsymbol{x}^{k+1} + \boldsymbol{u}^k \right]_{\mathcal{X}} \qquad\qquad k = 0, 1, 2, \dots$$
$$\boldsymbol{u}^{k+1} = \boldsymbol{u}^k + \left( \boldsymbol{x}^{k+1} - \boldsymbol{z}^{k+1} \right)$$

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

- **Explanation of steps:**
  - **$\boldsymbol{x}$-update:** Minimize $f(\boldsymbol{x})$ with a quadratic penalty term.
  - **$\boldsymbol{z}$-update:** Project $\boldsymbol{x}^{k+1} + \boldsymbol{u}^k$ onto the set $\mathcal{X}$.
  - **$\boldsymbol{u}$-update:** Update the scaled dual variable $\boldsymbol{u}$.

- **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:**

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

- **Reformulated problem for ADMM:**

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

# ADMM example: $\ell_2 - \ell_1$-norm minimization

- **ADMM algorithm:**
  - **$x$-update:**
    - Given $z$ and scaled dual variable $u$, solve:

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

    - Solution:

    $$x = \left(\boldsymbol{A}^\mathsf{T}\boldsymbol{A} + \rho\boldsymbol{I}\right)^{-1}\left(\boldsymbol{A}^\mathsf{T}\boldsymbol{b} + \rho(z - u)\right)$$

  - **$z$-update:**
    - Given $x$ and $u$, solve:

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

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

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

  - **$u$-update:**
    - Update the scaled dual variable:

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

# ADMM example: $\ell_2 - \ell_1$-norm minimization

- **ADMM iterative algorithm:**
  - Update rules:

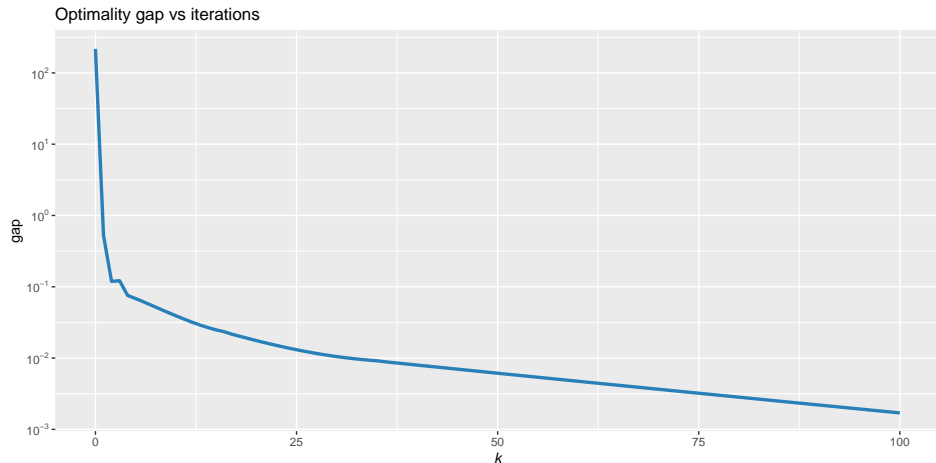$$\begin{aligned}
\mathbf{x}^{k+1} &= \left(\mathbf{A}^\mathsf{T}\mathbf{A} + \rho\mathbf{I}\right)^{-1}\left(\mathbf{A}^\mathsf{T}\mathbf{b} + \rho\left(\mathbf{z}^k - \mathbf{u}^k\right)\right) \\
\mathbf{z}^{k+1} &= \mathcal{S}_{\lambda/\rho}\left(\mathbf{x}^{k+1} + \mathbf{u}^k\right) \qquad\qquad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

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

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

# ADMM example: $\ell_2 - \ell_1$-norm minimization

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

# Outline

## Numerical comparison

- $\ell_2 - \ell_1$-**norm minimization problem:**

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad \tfrac{1}{2}\|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2 + \lambda\|\boldsymbol{x}\|_1.$$

- **Iterates of various algorithms:**
  - **BCD (Gauss-Seidel) iterates:**

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

  - $\boldsymbol{x}^{(k,i)} \triangleq \left(x_1^{k+1}, \ldots, x_{i-1}^{k+1}, x_i^k, \ldots, x_n^k\right)$

  - **Parallel BCD (Jacobi) iterates:**

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

## Numerical comparison

- **Iterates of various algorithms: (cont'd)**
  - **MM iterates:**

$$\boldsymbol{x}^{k+1} = \mathcal{S}_{\frac{\lambda}{\kappa}}\left(\boldsymbol{x}^k - \frac{1}{\kappa}\boldsymbol{A}^\mathsf{T}\left(\boldsymbol{A}\boldsymbol{x}^k - \boldsymbol{b}\right)\right), \qquad k = 0, 1, 2, \ldots$$

  - **Accelerated MM iterates:**

$$\begin{aligned}
\boldsymbol{r}^k &= R(\boldsymbol{x}^k) \triangleq \mathsf{MM}(\boldsymbol{x}^k) - \boldsymbol{x}^k \\
\boldsymbol{v}^k &= R(\mathsf{MM}(\boldsymbol{x}^k)) - R(\boldsymbol{x}^k) \\
\alpha^k &= -\max\left(1, \|\boldsymbol{r}^k\|_2/\|\boldsymbol{v}^k\|_2\right) \qquad k = 0, 1, 2, \ldots \\
\boldsymbol{y}^k &= \boldsymbol{x}^k - \alpha^k \boldsymbol{r}^k \\
\boldsymbol{x}^{k+1} &= \mathsf{MM}(\boldsymbol{y}^k)
\end{aligned}$$

## Numerical comparison

- **Iterates of various algorithms: (cont'd)**
  - **SCA iterates:**

$$\hat{x}^{k+1} = \mathcal{S}_{\frac{\lambda}{\tau + \text{diag}(A^\mathsf{T} A)}} \left( x^k - \frac{A^\mathsf{T} \left( A x^k - b \right)}{\tau + \text{diag}\left( A^\mathsf{T} A \right)} \right) \qquad k = 0, 1, 2, \ldots$$

$$x^{k+1} = \gamma^k \hat{x}^{k+1} + \left( 1 - \gamma^k \right) x^k$$

  - **ADMM iterates:**

$$x^{k+1} = \left( A^\mathsf{T} A + \rho I \right)^{-1} \left( A^\mathsf{T} b + \rho \left( z^k - u^k \right) \right)$$

$$z^{k+1} = \mathcal{S}_{\lambda/\rho} \left( x^{k+1} + u^k \right) \qquad k = 0, 1, 2, \ldots$$

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

## Numerical comparison

- **Comparison of methods:**
  - **BCD:**
    - Updates each element sequentially ($n = 100$).
    - High computational cost (CPU time) due to sequential updates.
  - **Jacobi:**
    - Parallel version of BCD.
    - Not guaranteed to converge.
    - Similar to SCA but lacks $\tau$ and smoothing step.
  - **MM:**
    - Requires computing the largest eigenvalue of $\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}$.
    - Conservative upper-bound $\kappa$ used for all elements.
  - **SCA:**
    - Uses $\mathrm{diag}\left(\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}\right)$ instead of a common $\kappa$.
    - Faster convergence due to element-specific updates.
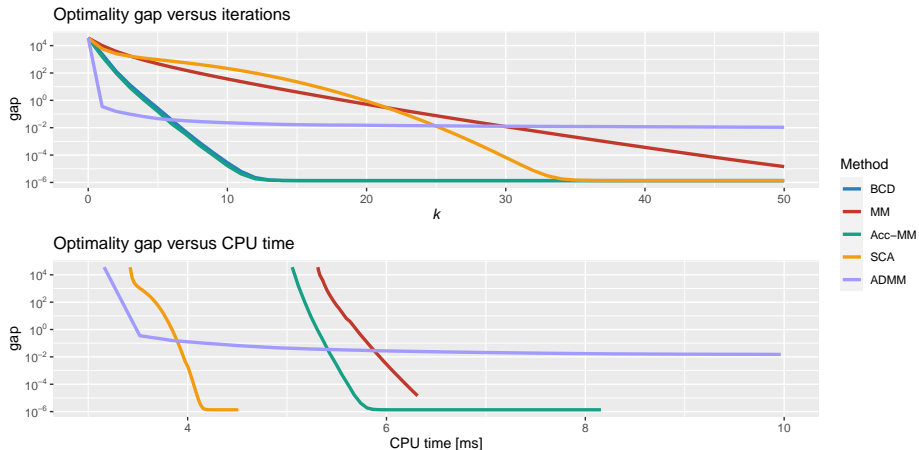  - **ADMM:**
    - Converges with lower accuracy.
    - Often sufficient for practical applications.

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



Optimality gap versus iterations

Optimality gap versus CPU time

Method
- BCD
- MM
- Acc–MM
- SCA
- ADMM

# Outline

# 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)

# References I

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.

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.

———. 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.

# References III

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. 2024. *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.