

Lecture 12: Linearly constrained
nonlinear optimization

Michael Patrinos

26 February 2004

nonlinear.

- Consider a constraint " $y_i(\mathbf{x}) \leq b_i$ ", where y_i is $\{x_k\}$ is infeasible until convergence. Why?
defining X ; in some cases even such that the sequence
- Most methods for (1) manipulate the constraints

is C_1 on X .

$X \subseteq \mathbb{R}^n$ non-empty, closed and convex set; $f : \mathbb{R}^n \rightarrow \mathbb{R}$

(1b) subject to $\mathbf{x} \in X$,

(1a) $\cdot (\mathbf{x}) = f_* = \min f$

- Consider the problem to find

Feasible-direction methods

- Checking whether \mathbf{p} is a feasible direction at x , or what the maximum feasible step from x in the direction of \mathbf{p} is, is very difficult.
- For which step length $a < 0$ does it happen that $g_i(\mathbf{x} + a\mathbf{p}) = b_i$? This is a nonlinear equation in a !
- Assuming that X is Polyhedral, these problems are not present.

Feasible-direction descent methods

Step 0. Determine a starting point $\mathbf{x}^0 \in \mathbb{R}^n$ such that

$$\mathbf{x}^0 \in X. \text{ Set } k := 0.$$

Step 1. Determine a search direction $\mathbf{p}^k \in \mathbb{R}^n$ such that \mathbf{p}^k is a feasible direction.

$$f(\mathbf{x}^k + \alpha_k \mathbf{p}^k) > f(\mathbf{x}^k) \text{ and } \mathbf{x}^k + \alpha_k \mathbf{p}^k \in X.$$

Step 2. Determine a step length $\alpha_k < 0$ such that

Step 3. Let $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{p}^k$.

Step 4. If a termination criterion is fulfilled, then stop!

Otherwise, let $k := k + 1$ and go to Step 1.

- Similar form as the general method for unconstrained optimization.
- Just as *Local* as methods for unconstrained optimization.
- Search directions typically based on the approximation of f —relaxation.
- Line searches similar; note the maximum step.
- Termination criteria and descent based on first-order optimality (remember the unconstrained condition that $\Delta f(\mathbf{x}_*) = \mathbf{0}$ holds).

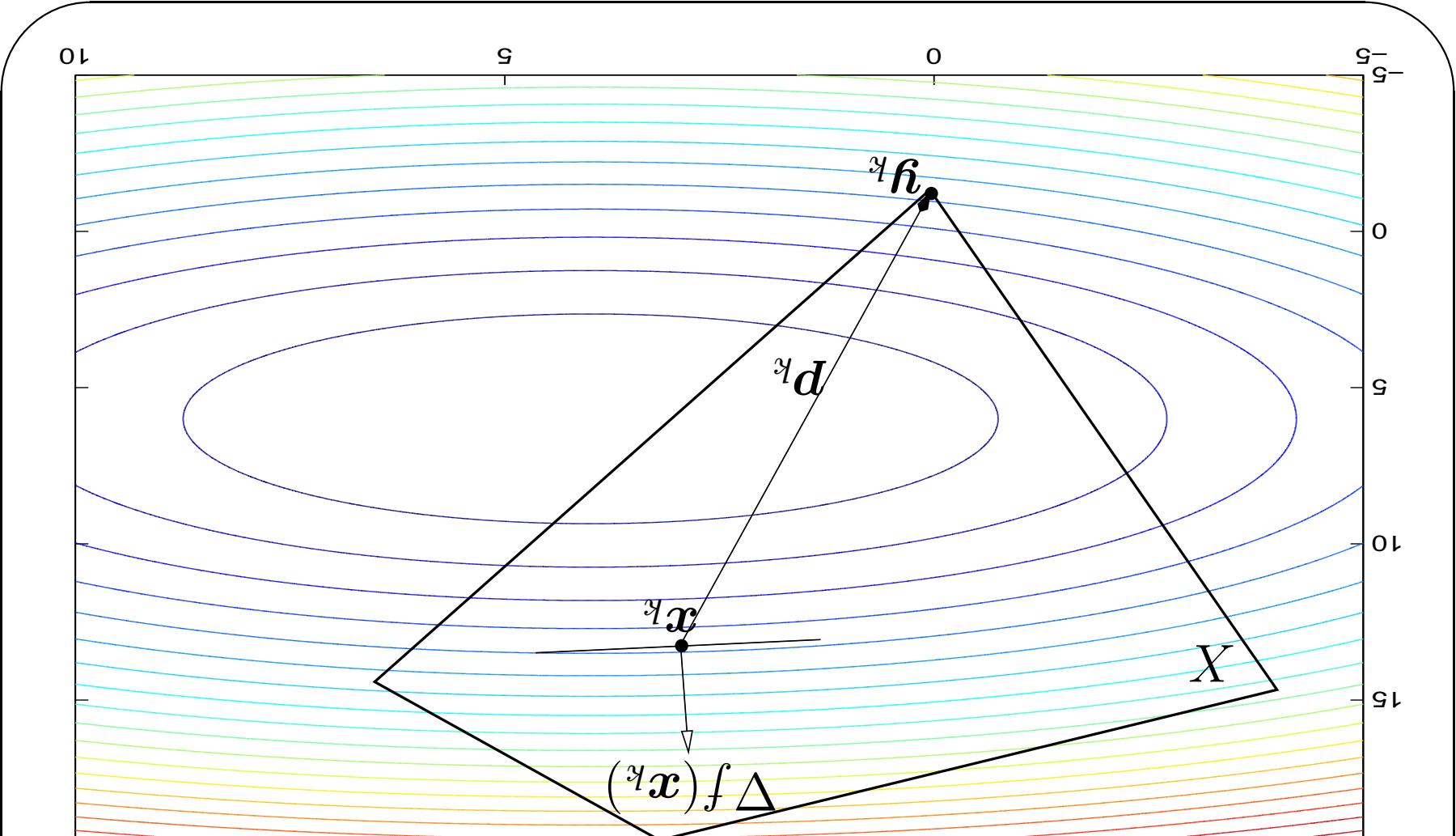
Notes

LP-based algorithm, I: The Frank-Wolfe method

- The Frank-Wolfe (1952) method is based on a first-order approximation of f around the iterate \mathbf{x}^k . This means that the relaxed problems are LPs, which can then be solved by using the Simplex method.
- Remember the following first-order condition [Proposition 4.20(b)]: If $\mathbf{x}_* \in X$ is a local minimum of f on X then
$$(\mathcal{L}_{\mathbf{x}}(\mathbf{x}_*))' = 0,$$
 holds.

- This is the basis of the Frank-Wolfe algorithm.
- The search direction is towards an extreme point [one that is optimal in the LP over X with costs $\Delta f(\mathbf{x}^k)$].
- The search direction is given by $\mathbf{p}_k := \mathbf{y}_k - \mathbf{x}_k$: \mathbf{y}_k is a feasible descent direction and \mathbf{y}_k is a solution to this LP problem, then the direction of \mathbf{p}_k is with respect to f at \mathbf{x} .
- Follows that if, given an iterate $\mathbf{x}^k \in X$,
- (3)
$$0 = (\mathbf{x}^* - \mathbf{x}^k)^\top (\mathbf{x}^* - \mathbf{x}^k) \geq \min_{\mathbf{x} \in X} f(\mathbf{x})$$
- Remember also the following equivalent statement:

- Note: We must assume that X is bounded in order to ensure that the LP always has a finite solution. The algorithm can in fact be extended to allow for unbounded solutions to the LP, and thereby extending the Frank-Wolfe method for general polyhedra; the search directions then are either towards an extreme point (finite solution to LP) or in the direction of an extreme ray of X (unbounded solution to LP).



The search-direction problem

terminate! Otherwise, let $k := k + 1$ and go to Step 1.

Step 4. If, for example, $z^k(\mathbf{y}^k)$ or a^k is close to zero, then

Step 3. Let $\mathbf{x}^{k+1} := \mathbf{x}^k + a^k \mathbf{d}^k$.

$f(\mathbf{x}^k + a^k \mathbf{d}^k)$ over $a \in [0, 1]$. Let a^k be the step length.

Step 2. Approximately solve the problem to minimize

Let $\mathbf{d}^k := \mathbf{y}^k - \mathbf{x}^k$ be the search direction.

$$(4) \quad \text{minimize}_{\mathbf{y} \in X} z^k(\mathbf{y}) =: (\mathbf{y}^k)^T \Delta (\mathbf{y}^k - \mathbf{x}^k).$$

Step 1. Find a solution \mathbf{y}^k to the problem to

X . Set $k := 0$.

Step 0. Find $\mathbf{x}^0 \in X$, for example any extreme point in

Algorithm description, Frank-Wolfe

- Theorem 12.1: Suppose that $X \subseteq \mathbb{R}^n$ is a non-empty, bounded polyhedron, and that the function f is in C_1 on X . Suppose that in Step 2 of the Frank-Wolfe algorithm, we either use an exact line search or the Armijo step length rule. Then, the sequence $\{\mathbf{x}_k\}$ is bounded, $\{f(\mathbf{x}_k)\}$ is descending, and every limit point (at least one exists) is stationary; further, the sequence $\{z_k(\mathbf{y}_k)\} \rightarrow 0$.
- If f is convex on X , then every limit point is globally optimal.

Convergence

The convex case: Lower bounds

- Remember the following characterization of convex functions in C_1 on X [Theorem 3.44(a)]:

f is convex on $X \iff$

$$\cdot X \ni \mathbf{y}, \mathbf{x} \text{ for all } \mathbf{x}, \mathbf{y} \in X. \quad f(\mathbf{x}) + f(\mathbf{y}) \geq f(\mathbf{x} + \mathbf{y})$$

- Suppose f is convex on X . Then, $f(\mathbf{x}) \geq f(\mathbf{z}) + (\nabla f(\mathbf{z}))^\top (\mathbf{x} - \mathbf{z})$ if and only if \mathbf{z} is globally optimal.

- Utilize the lower bound as follows: we know that check in Step 4 whether $|f(\mathbf{x}_k) - \text{LBD}| / |\text{LBD}|$ is small, store the best LBD, and $f_* \in [f(\mathbf{x}_k) + (\nabla f(\mathbf{x}_k))^\top \mathbf{z}_k, f(\mathbf{x}_k)]$. Store the best LBD, and and if so terminate.

- Another reason is that the information generated (the linear subproblems, we can do better by storing and utilizing this information.
- For highly nonlinear problems, the approximation is bad—the optimal solution may be far from an extreme point.
- In order to find a near-optimum requires many iterations—the algorithm is slow.
- Frank-Wolfe uses linear approximations—works best for almost linear problems.

Final comments



for some $\alpha_1, \dots, \alpha_k \geq 0$ such that $\sum_{i=1}^k \alpha_i = 1$.

$$\sum_{i=1}^k \alpha_i \mathbf{a}_i = \mathbf{x}$$

- convex combination of the points in V , that is,
points of P . Every $\mathbf{x} \in P$ can be represented as a
bounded, and $V = \{\mathbf{a}_1, \dots, \mathbf{a}_K\}$ be the set of extreme
 $P = \{\mathbf{x} \in \mathbb{R}^n \mid A\mathbf{x} = \mathbf{q}; \mathbf{x} \geq \mathbf{0}_n\}$, be non-empty and
case for bounded polyhedra): Let
- Remember the Representation Theorem 9.9 (special

LP-based algorithm, II: Simplicial decomposition

- The process is still iterative: we generate a „working set“ P_k of indices i , optimize the function f over the convex hull of the known points, and check for stationarity and/or generate a new extreme point.
- The idea behind the Simplex decomposition method is to generate the extreme points \mathbf{u}_i which can be used to describe an optimal solution \mathbf{x}_* , that is, the vectors \mathbf{u}_i with positive weights a_i in
$$\sum_{i=1}^K a_i \mathbf{u}_i = \mathbf{x}_*$$

Algorithm description, Simplicial decomposition

Step 0. Find $x^0 \in X$, for example any extreme point in

X . Set $k := 0$. Let $P_0 := \emptyset$.

Step 1. Let y^k be a solution to the LP problem (4).
Let $P_{k+1} := P_k \cup \{k\}$.

and go to Step 1.

$P^{k+1} = P^k$, then terminate! Otherwise, let $k := k + 1$

Step 4. If, for example, $z^k(\mathbf{y}^k)$ is close to zero, or if

Step 3. Let $\mathbf{x}^{k+1} := \mathbf{u}^{k+1}\mathbf{x}^k + \sum_{i \in P^{k+1}} \mathbf{u}^{k+1} \mathbf{y}_i$.

$$(5c) \quad u_i < 0, \quad i \in P^{k+1}.$$

$$(5d) \quad \text{subject to } u_i = 1, \quad \sum_{i \in P^{k+1}} u_i + \sum_{i \in P^{k+1}} u_i = 1,$$

$$(5e) \quad \text{minimize}_{(\mathbf{u}, \boldsymbol{\alpha})} \left(\mathbf{u}^T \sum_{i \in P^{k+1}} \mathbf{y}_i + \boldsymbol{\alpha}^T \mathbf{x}^k \right) \text{ s.t.}$$

restricted master problem to

Step 2. Let $(\mathbf{u}^k, \boldsymbol{\alpha}^{k+1})$ be an approximate solution to the

- This basic algorithm keeps all information generated, and adds one new extreme point in every iteration.
- An alternative is to drop columns (vectors y_i) that have received a zero weight, or to keep only a maximum number of vectors. (Stated in the Notes.)
- Special case: maximum number of vectors kept = 1
 ⇔ the Frank-Wolfe algorithm
- We obviously improve the Frank-Wolfe algorithm by utilizing more information.

Convergence

- Based on the fact that it does at least as well as the Frank-Wolfe algorithm.

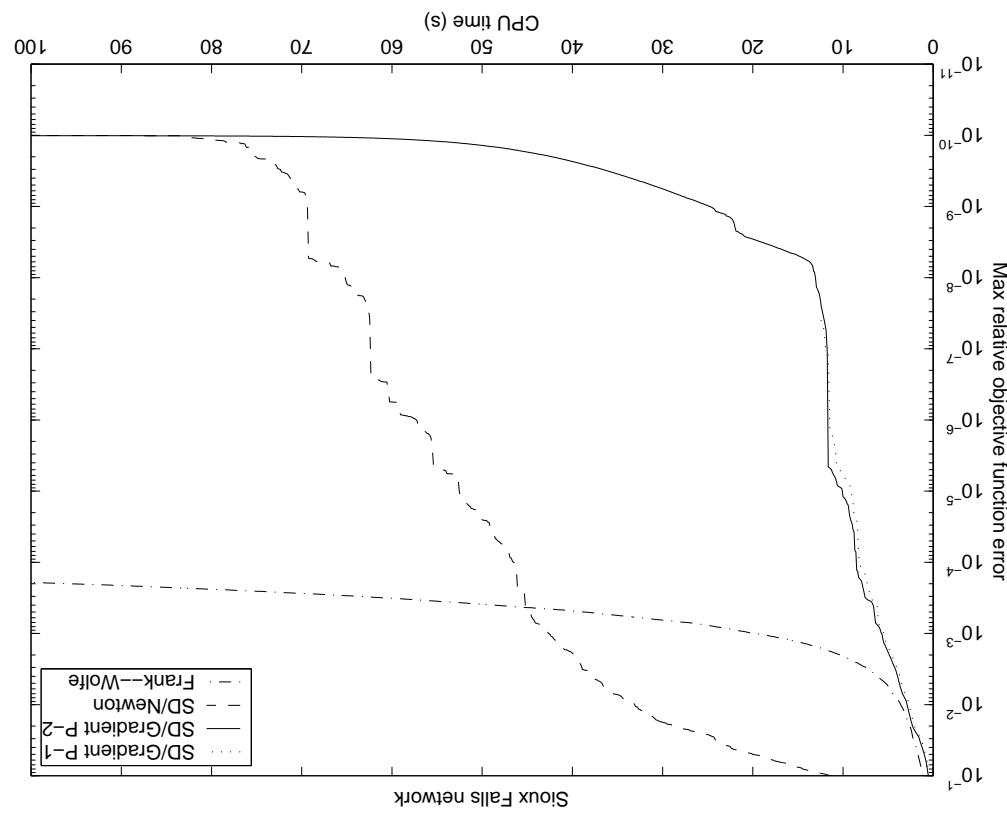
Convergence is finite if the restricted master problems (RMPs) are solved exactly, and the maximum number of vectors kept is at least as many as are needed to span x_* .

- Much more efficient than the Frank-Wolfe algorithm in practice.
- We can solve the RMPs efficiently, since they are almost unconstrained.

- A large-scale non-linear network flow problem which is used to estimate traffic flows in cities.
- The model is over the small city of Sioux Falls in North Dakota, whose representation has 24 nodes, 76 links, and 528 pairs of origin and destination.
- Three algorithms for the RMPs were tested—a Newton method and two gradient projection methods (see the next section). A MATLAB implementation.
- Remarkable difference—The Frank-Wolfe method suffers from very small step length being taken.

An illustration of FW vs. SD

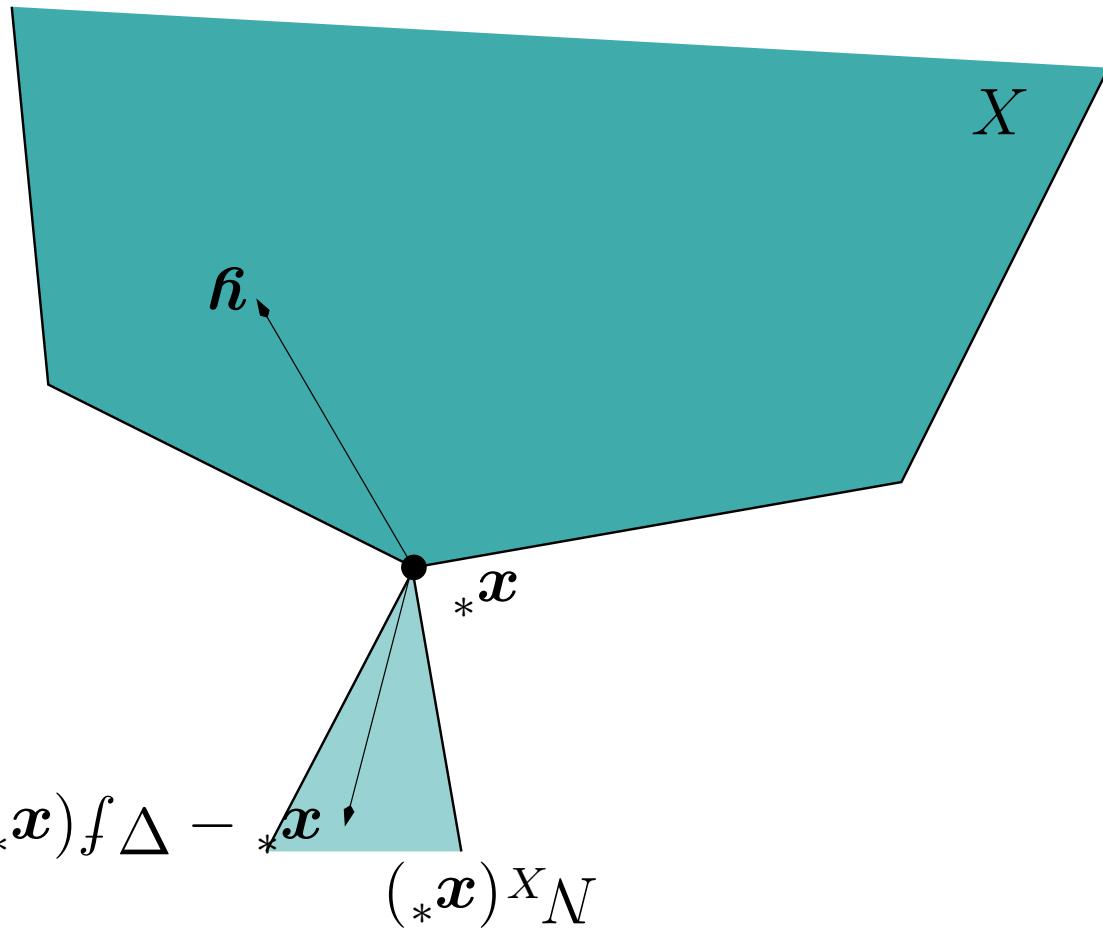
Figure 1: The performance of DSD vs. FW on the Sioux Falls network.



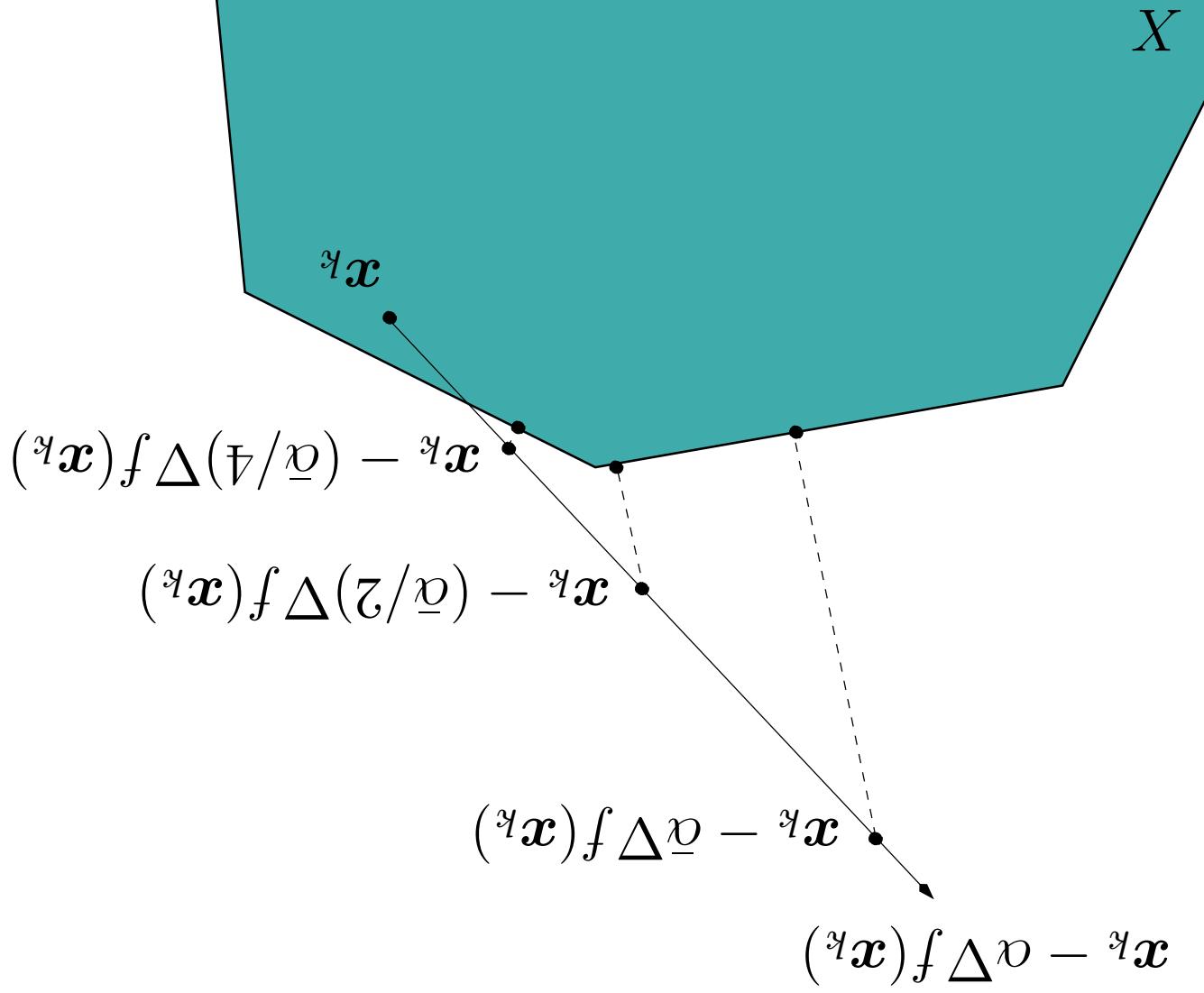
- The gradient projection algorithm is based on the projection characterization of a stationary point (Section 4.4): \mathbf{x}_* is a stationary point if and only if

$$\cdot [(_*\mathbf{x}) f \Delta - {}_*\mathbf{x}]^X = {}_*\mathbf{x}$$

QP-based algorithm: The gradient projection algorithm



- Let $\mathbf{d} := \text{Proj}_X[\mathbf{x} - \alpha \Delta f(\mathbf{x})] - \mathbf{x}$, for any $\alpha < 0$. Then, if and only if \mathbf{x} is non-stationary, \mathbf{d} is a feasible descent direction of f at \mathbf{x} .
 - The gradient projection algorithm is normally stated such that the line search is done over the projection arc, that is, we find a step length α_k for which has a good objective value. Use the Armijo rule to determine α_k :
- $$(6) \quad \mathbf{x}_{k+1} := \text{Proj}_X[\mathbf{x}_k - \alpha_k \Delta f(\mathbf{x}_k)], \quad k = 1, \dots.$$



- Theorem 12.3 (simplified): Suppose that $X \subseteq \mathbb{R}^n$ is non-empty, compact and convex. Consider the iterative algorithm defined by the iteration (6), where the step length α_k is determined by the Armijo step length rule along the projection arc. Then, the sequence $\{\mathbf{x}_k\}$ is bounded, the sequence $\{f(\mathbf{x}_k)\}$ is descending, lower bounded and therefore has a limit, and every limit point of $\{\mathbf{x}_k\}$ is stationary.
- Gradient projection becomes steepest descent with Armijo line search when $X = \mathbb{R}^n$.
- Convergence arguments similar to steepest descent one.

Convergence

- Set up a Phase I problem for the linear inequality system, and treat the complementarity conditions implicitly, as follows: if x_j (respectively, u_j) is a basic variable, then u_j (respectively, x_j) must not be an entering variable. Same for the pair (s_i, u_i) .

- Result: A system of linear inequalities and equalities plus two sets of complementarity conditions of the form $x_j u_j = 0$ and $s_i u_i = 0$.
- Add slack variables.

- State the KKT conditions for the strictly convex QP problem which determines the projection.

Quadratic subproblems—how are they solved?