

Algorithms for constrained optimization

- Projected Gradient Methods with Linear Constraints

$$\begin{aligned} & \min \quad f(\mathbf{x}) \\ & s.t. \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned} \tag{1}$$

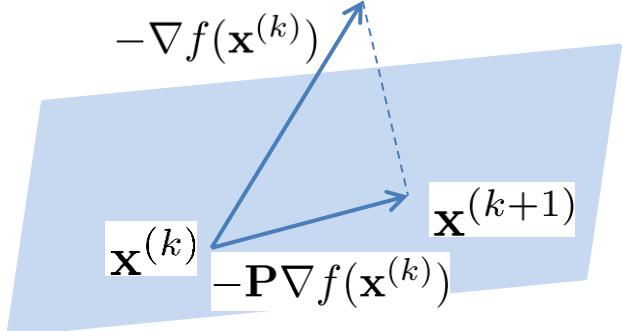
where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m < n$, $\text{rank } \mathbf{A} = m$, $\mathbf{b} \in \mathbb{R}^m$.

$$\mathbf{P} = \mathbf{I}_n - \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{A},$$

$$\mathcal{P}_\Omega(\mathbf{x}) = \mathbf{x} - \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}(\mathbf{A}\mathbf{x} - \mathbf{b}) = \mathbf{P}\mathbf{x} + \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{b}.$$

$$\mathbf{x}^{(k)} \in \Omega \implies \mathbf{x}^{(k+1)} = \mathcal{P}_\Omega(\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)}))$$

$$\begin{aligned} &= [\mathbf{I} - \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{A}] (\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})) \\ &\quad + \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{b} \\ &= \mathbf{x}^{(k)} - \alpha_k \mathbf{P} \nabla f(\mathbf{x}^{(k)}). \end{aligned}$$



Algorithms for constrained optimization

- Projected Gradient Methods with Linear Constraints

Proposition 1. $-\mathbf{P}\nabla f(\mathbf{x}^{(k)})$ points in the direction of maximum rate of decrease of f at $\mathbf{x}^{(k)}$ along the surface defined by $\mathbf{Ax} = \mathbf{b}$.

Proposition 2. If $-\mathbf{P}\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ then it is a descent direction.

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

$$\begin{aligned} \min_{\mathbf{x}} \quad & \frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x} + \mathbf{q}^T \mathbf{x} + r, \\ \text{s.t. } \quad & \mathbf{A} \mathbf{x} = \mathbf{b}, \end{aligned}$$

where $\mathbf{P} \in \mathbb{S}_+^n$. The optimality conditions are:

$$\mathbf{A} \mathbf{x}^* = \mathbf{b}, \quad \mathbf{P} \mathbf{x}^* + \mathbf{q} + \mathbf{A}^T \boldsymbol{\nu}^* = \mathbf{0},$$

which can be written as

$$\boxed{\begin{matrix} \text{KKT matrix} & \xrightarrow{\hspace{1cm}} & \boxed{\begin{pmatrix} \mathbf{P} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{x}^* \\ \boldsymbol{\nu}^* \end{pmatrix} = \begin{pmatrix} -\mathbf{q} \\ \mathbf{b} \end{pmatrix}} & \xrightarrow{\hspace{1cm}} \text{KKT system} \end{matrix}}$$

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

Nonsingularity of the KKT matrix

Recall our assumption that $\mathbf{P} \in \mathbb{S}_+^n$ and $\text{rank } \mathbf{A} = p < n$. There are several conditions equivalent to nonsingularity of the KKT matrix:

- $\mathcal{N}(\mathbf{P}) \cap \mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$, i.e., \mathbf{P} and \mathbf{A} have no nontrivial common nullspace.
- $\mathbf{Ax} = \mathbf{0}, \mathbf{x} \neq \mathbf{0} \implies \mathbf{x}^T \mathbf{Px} > 0$, i.e., \mathbf{P} is positive definite on the nullspace of \mathbf{A} .
- $\mathbf{F}^T \mathbf{PF} \succ \mathbf{0}$, where $\mathbf{F} \in \mathbb{R}^{n \times (n-p)}$ is a matrix for which $\mathcal{R}(\mathbf{F}) = \mathcal{N}(\mathbf{A})$.

As an important special case, we note that if $\mathbf{P} \succ \mathbf{0}$, the KKT matrix must be nonsingular.

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

Solving KKT systems

$$\begin{pmatrix} \mathbf{H} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix} = - \begin{pmatrix} \mathbf{g} \\ \mathbf{h} \end{pmatrix}. \quad (2)$$

Here we assume $\mathbf{H} \in \mathbb{S}_+^n$ and $\mathbf{A} \in \mathbb{R}^{p \times n}$ with $\text{rank } \mathbf{A} = p < n$.

Solving KKT systems

One straightforward approach is to simply solve the KKT system (1), which is a set of $n+p$ linear equations in $n+p$ variables. The KKT matrix is symmetric, but may not be positive definite, so a good way to do this is to use an LDL^T factorization. If no structure of the matrix is exploited, the cost is $(1/3)(n+p)^3$ flops. This can be a reasonable approach when the problem is small (i.e., n and p are not too large), or when \mathbf{A} and \mathbf{H} are sparse.

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

Solving KKT system via elimination

We start by describing the simplest case, in which $\mathbf{H} \succ \mathbf{0}$. Starting from the first of the KKT equations

$$\mathbf{H}\mathbf{v} + \mathbf{A}^T\mathbf{w} = -\mathbf{g}, \quad \mathbf{A}\mathbf{v} = -\mathbf{h},$$

we solve for \mathbf{v} to obtain

$$\mathbf{v} = -\mathbf{H}^{-1}(\mathbf{g} + \mathbf{A}^T\mathbf{w}).$$

Substituting this into the second KKT equation yields $\mathbf{A}\mathbf{H}^{-1}(\mathbf{g} + \mathbf{A}^T\mathbf{w}) = \mathbf{h}$, so we have

$$\mathbf{w} = (\mathbf{A}\mathbf{H}^{-1}\mathbf{A}^T)^{-1}(\mathbf{h} - \mathbf{A}\mathbf{H}^{-1}\mathbf{g}).$$

These formulae give us a method for computing \mathbf{v} and \mathbf{w} .

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

Solving KKT system by block elimination.

Given: KKT system with $\mathbf{H} \succ 0$.

1. Form $\mathbf{H}^{-1}\mathbf{A}^T$ and $\mathbf{H}^{-1}\mathbf{g}$.
2. Form Schur complement $\mathbf{S} = -\mathbf{A}\mathbf{H}^{-1}\mathbf{A}^T$.
3. Determine \mathbf{w} by solving $\mathbf{Sw} = \mathbf{AH}^{-1}\mathbf{g} - \mathbf{h}$.
4. Determine \mathbf{v} by solving $\mathbf{Hv} = -\mathbf{A}^T\mathbf{w} - \mathbf{g}$.

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

Step 1 can be done by a Cholesky factorization of \mathbf{H} , followed by $p+1$ solves, which costs $f + (p+1)s$, where f is the cost of factoring \mathbf{H} and s is the cost of an associated solve. Step 2 requires a $p \times n$ by $n \times p$ matrix multiplication. If we exploit no structure in this calculation, the cost is p^2n flops. (Since the result is symmetric, we only need to compute the upper triangular part of \mathbf{S} .) In some cases special structure in \mathbf{A} and \mathbf{H} can be exploited to carry out step 2 more efficiently. Step 3 can be carried out by Cholesky factorization of $-\mathbf{S}$, which costs $(1/3)p^3$ flops if no further structure of \mathbf{S} is exploited. Step 4 can be carried out using the factorization of \mathbf{H} already calculated in step 1, so the cost is $2np + s$ flops. The total flop count, assuming that no structure is exploited in forming or factoring the Schur complement, is

$$f + ps + p^2n + (1/3)p^3$$

flops (keeping only dominant terms). If we exploit structure in forming or factoring \mathbf{S} , the last two terms are even smaller.

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

If \mathbf{H} can be factored efficiently, then block elimination gives us a flop count advantage over directly solving the KKT system using an LDL^T factorization. For example, if \mathbf{H} is diagonal (which corresponds to a separable objective function), we have $f = 0$ and $s = n$, so the total cost is $p^2n + (1/3)p^3$ flops, which grows only linearly with n . If \mathbf{H} is banded with bandwidth $k \ll n$, then $f = nk^2$, $s = 4nk$, so the total cost is around $nk^2 + 4nkp + p^2n + (1/3)p^3$ which still grows only linearly with n . Other structures of \mathbf{H} that can be exploited are block diagonal (which corresponds to block separable objective function), sparse, or diagonal plus low rank.

Examples: Equality constrained analytic center, Minimum length piecewise-linear curve subject to equality constraints, Locally linear embedding (LLE)

Algorithms for constrained optimization

- Equality constrained convex quadratic minimization

Elimination with singular \mathbf{H}

The block elimination method described above obviously does not work when \mathbf{H} is singular, but a simple variation on the method can be used in this more general case. The more general method is based on the following result: The KKT matrix is nonsingular iff $\mathbf{H} + \mathbf{A}^T \mathbf{Q} \mathbf{A} \succ \mathbf{0}$ for some $\mathbf{Q} \succeq \mathbf{0}$, in which case, $\mathbf{H} + \mathbf{A}^T \mathbf{Q} \mathbf{A} \succ \mathbf{0}$ for all $\mathbf{Q} \succ \mathbf{0}$. We conclude, for example, that if the KKT matrix is nonsingular, then $\mathbf{H} + \mathbf{A}^T \mathbf{A} \succ \mathbf{0}$.

Let $\mathbf{Q} \succeq \mathbf{0}$ be a matrix for which $\mathbf{H} + \mathbf{A}^T \mathbf{Q} \mathbf{A} \succ \mathbf{0}$. Then the KKT system (2) is equivalent to

$$\begin{pmatrix} \mathbf{H} + \mathbf{A}^T \mathbf{Q} \mathbf{A} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix} = - \begin{pmatrix} \mathbf{g} + \mathbf{A}^T \mathbf{Q} \mathbf{h} \\ \mathbf{h} \end{pmatrix}, \quad (3)$$

which can be solved using elimination since $\mathbf{H} + \mathbf{A}^T \mathbf{Q} \mathbf{A} \succ \mathbf{0}$.

Examples: Equality constrained analytic centering, Optimal network flow, Optimal network flow, Analytic center of a linear matrix inequality

Algorithms for constrained optimization

- Eliminating equality constraints

One general approach to solving the equality constrained problem is to eliminate the equality constraints and then solve the resulting unconstrained problem using methods for unconstrained minimization. We first find a matrix $\mathbf{F} \in \mathbb{R}^{n \times (n-p)}$ and vector $\hat{\mathbf{x}} \in \mathbb{R}^n$ that parameterize the (affine) feasible set:

$$\{\mathbf{x} | \mathbf{A}\mathbf{x} = \mathbf{b}\} = \{\mathbf{F}\mathbf{z} + \hat{\mathbf{x}} | \mathbf{z} \in \mathbb{R}^{n-p}\}.$$

Here $\hat{\mathbf{x}}$ can be chosen as any particular solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $\mathbf{F} \in \mathbb{R}^{n \times (n-p)}$ is any matrix whose range is the nullspace of \mathbf{A} . We then form the reduced or eliminated optimization problem:

$$\min_{\mathbf{z}} \tilde{f}(\mathbf{z}) \triangleq f(\mathbf{F}\mathbf{z} + \hat{\mathbf{x}}), \quad (4)$$

which is an unconstrained problem with variable $\mathbf{z} \in \mathbb{R}^{n-p}$. From its solution \mathbf{z}^* , we can find the solution of the equality constrained problem as $\mathbf{x}^* = \mathbf{F}\mathbf{z}^* + \hat{\mathbf{x}}$.

Algorithms for constrained optimization

- Eliminating equality constraints

Example: Optimal allocation with resource constraint. We consider the problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \sum_{i=1}^n f_i(x_i), \\ \text{s.t.} \quad & \sum_{i=1}^n x_i = b, \end{aligned}$$

where the functions $f_i : \mathbb{R} \rightarrow \mathbb{R}$ are convex and twice differentiable, and $b \in \mathbb{R}$ is a problem parameter. We interpret this as the problem of optimally allocating a single resource, with a fixed total amount b (the budget) to n otherwise independent activities.

Algorithms for constrained optimization

- Solving equality constrained problems via the dual

Another approach to solving (1) is to solve the dual, and then recover the optimal primal variable \mathbf{x}^* . The dual function of (1) is:

$$\begin{aligned} g(\boldsymbol{\nu}) &= -\mathbf{b}^T \boldsymbol{\nu} + \inf_{\mathbf{x}} (f(\mathbf{x}) + \boldsymbol{\nu}^T \mathbf{A}\mathbf{x}) \\ &= -\mathbf{b}^T \boldsymbol{\nu} - \sup_{\mathbf{x}} ((-\mathbf{A}^T \boldsymbol{\nu})^T \mathbf{x} - f(\mathbf{x})) = -\mathbf{b}^T \boldsymbol{\nu} - f^*(-\mathbf{A}^T \boldsymbol{\nu}), \end{aligned}$$

where f^* is the conjugate of f . So the dual problem is:

$$\max_{\boldsymbol{\nu}} -\mathbf{b}^T \boldsymbol{\nu} - f^*(-\mathbf{A}^T \boldsymbol{\nu}).$$

Since by assumption there is an optimal point, the problem is strictly feasible, so Slater's condition holds. Therefore strong duality holds, and the dual optimum is attained, i.e., there exists a $\boldsymbol{\nu}^*$ with $g(\boldsymbol{\nu}^*) = p^*$.

Once we find an optimal dual variable $\boldsymbol{\nu}^*$, we reconstruct an optimal primal solution \mathbf{x}^* from it.

Algorithms for constrained optimization

- Solving equality constrained problems via the dual

Example: Equality constrained analytic center. We consider the problem

$$\min_{\mathbf{x}} f(\mathbf{x}) = - \sum_{i=1}^n \log x_i, \quad s.t. \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \quad (4)$$

where $\mathbf{A} \in \mathbb{R}^{p \times n}$, with implicit constraint $\mathbf{x} > \mathbf{0}$. Using

$$f^*(\mathbf{y}) = \sum_{i=1}^n (-1 - \log(-y_i)) = -n - \sum_{i=1}^n \log(-y_i)$$

(with $\text{dom } f^* = -\mathbb{R}_{++}^n$), the dual problem is

$$\max_{\boldsymbol{\nu}} g(\boldsymbol{\nu}) = -\mathbf{b}^T \boldsymbol{\nu} + n + \sum_{i=1}^n \log(\mathbf{A}^T \boldsymbol{\nu})_i, \quad (5)$$

with implicit constraint $\mathbf{A}^T \boldsymbol{\nu} > \mathbf{0}$.

Algorithms for constrained optimization

- Solving equality constrained problems via the dual

Here we can easily solve the dual feasibility equation, i.e., find the \mathbf{x} that minimizes $L(\mathbf{x}, \boldsymbol{\nu})$:

$$\nabla f(\mathbf{x}) + \mathbf{A}^T \boldsymbol{\nu} = -(1/x_1, \dots, 1/x_n) + \mathbf{A}^T \boldsymbol{\nu} = \mathbf{0},$$

and so

$$x_i(\boldsymbol{\nu}) = 1/(\mathbf{A}^T \boldsymbol{\nu})_i. \tag{6}$$

To solve the equality constrained analytic centering problem (4), we solve the (unconstrained) dual problem (5), and then recover the optimal solution of (4) via (6).

Algorithms for constrained optimization

- Newton's method with equality constraints

We describe an extension of Newton's method to include equality constraints. The method is almost the same as Newton's method without constraints, except for two differences: The initial point must be feasible (i.e., satisfy $\mathbf{x} \in \text{dom } f$ and $\mathbf{Ax} = \mathbf{b}$), and the definition of Newton step is modified to take the equality constraints into account. In particular, we make sure that the Newton step $\Delta\mathbf{x}_{nt}$ is a feasible direction, i.e., $\mathbf{A}\Delta\mathbf{x}_{nt} = \mathbf{0}$.

To derive the Newton step $\Delta\mathbf{x}_{nt}$ for the equality constrained problem (1) at the feasible point \mathbf{x} , we replace the objective with its second-order Taylor approximation near \mathbf{x} , to form the problem

$$\begin{aligned} \min_{\mathbf{v}} \hat{f}(\mathbf{x} + \mathbf{v}) &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v} + (1/2)\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v}, \\ s.t. \quad \mathbf{A}(\mathbf{x} + \mathbf{v}) &= \mathbf{b}. \end{aligned} \tag{7}$$

Algorithms for constrained optimization

- Newton's method with equality constraints

This is a (convex) quadratic minimization problem with equality constraints, and can be solved analytically. We define $\Delta\mathbf{x}_{nt}$, the Newton step at \mathbf{x} , as the solution of the convex quadratic problem (4), assuming the associated KKT matrix is nonsingular. In other words, the Newton step $\Delta\mathbf{x}_{nt}$ is what must be added to \mathbf{x} to solve the problem when the quadratic approximation is used in place of f .

From our previous analysis on linearly constrained quadratic problems, the Newton step $\Delta\mathbf{x}_{nt}$ is characterized by

$$\begin{pmatrix} \nabla^2 f(\mathbf{x}) & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \Delta\mathbf{x}_{nt} \\ \mathbf{w} \end{pmatrix} = \begin{pmatrix} -\nabla f(\mathbf{x}) \\ \mathbf{0} \end{pmatrix}, \quad (8)$$

where \mathbf{w} is the associated optimal dual variable for the quadratic problem. The Newton step is defined only at points for which the KKT matrix is nonsingular.

Algorithms for constrained optimization

- Newton's method with equality constraints

Solution of linearized optimality conditions: We can interpret the Newton step $\Delta\mathbf{x}_{nt}$, and the associated vector \mathbf{w} , as the solutions of a linearized approximation of the optimality conditions

$$\mathbf{A}\mathbf{x}^* = \mathbf{b}, \quad \nabla f(\mathbf{x}^*) + \mathbf{A}^T \boldsymbol{\nu}^* = \mathbf{0}.$$

We substitute $\mathbf{x} + \Delta\mathbf{x}_{nt}$ for \mathbf{x}^* and \mathbf{w} for $\boldsymbol{\nu}^*$, and replace the gradient term in the second equation by its linearized approximation near \mathbf{x} , to obtain the equations

$$\mathbf{A}(\mathbf{x} + \Delta\mathbf{x}_{nt}) = \mathbf{b}, \quad \nabla f(\mathbf{x} + \Delta\mathbf{x}_{nt}) + \mathbf{A}^T \mathbf{w} \approx \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x}) \Delta\mathbf{x}_{nt} + \mathbf{A}^T \mathbf{w} = \mathbf{0}.$$

Using $\mathbf{A}\mathbf{x} = \mathbf{b}$, these become

$$\mathbf{A}\Delta\mathbf{x}_{nt} = \mathbf{0}, \quad \nabla^2 f(\mathbf{x}) \Delta\mathbf{x}_{nt} + \mathbf{A}^T \mathbf{w} = -\nabla f(\mathbf{x}), \tag{9}$$

which are precisely the equations (8) that define the Newton step.

Algorithms for constrained optimization

- Newton's method with equality constraints

The Newton decrement:

Let

$$\lambda(\mathbf{x}) = (\Delta\mathbf{x}_{nt} \nabla^2 f(\mathbf{x}) \Delta\mathbf{x}_{nt})^{1/2}.$$

Then

$$f(\mathbf{x}) - \inf_{\mathbf{v}} \{\hat{f}(\mathbf{x} + \mathbf{v}) \mid \mathbf{A}(\mathbf{x} + \mathbf{v}) = \mathbf{b}\} = \frac{1}{2} \lambda^2(\mathbf{x}).$$

So $\frac{1}{2} \lambda^2(\mathbf{x})$ gives an estimate of $f(\mathbf{x}) - p^*$ and can serve as a good stopping criterion.

$\Delta\mathbf{x}_{nt}$ also gives a descent direction because by (9)

$$\nabla f(\mathbf{x})^T \Delta\mathbf{x}_{nt} = -\lambda^2(\mathbf{x}) < 0.$$

Algorithms for constrained optimization

- Newton's method with equality constraints

Newton's method for equality constrained minimization.

Given: starting point $\mathbf{x} \in \text{dom } f$ with $\mathbf{Ax} = \mathbf{b}$, tolerance $\epsilon > 0$.

Repeat:

1. Compute the Newton step and decrement $\Delta\mathbf{x}_{nt}$, $\lambda(\mathbf{x})$.
2. *Stopping criterion.* quit if $\lambda^2/2 \leq \epsilon$.
3. *Line search.* Choose step size t by backtracking line search.
4. *Update.* $\mathbf{x} := \mathbf{x} + t\Delta\mathbf{x}_{nt}$.

Algorithms for constrained optimization

- Dual Derivatives and Subgradients – Model Problem

We focus on the primal problem:

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & s.t. \mathbf{x} \in \mathcal{X}, g_j(\mathbf{x}) \leq 0, j = 1, \dots, r, \end{aligned} \tag{10}$$

and its dual

$$\begin{aligned} & \max_{\boldsymbol{\mu}} q(\boldsymbol{\mu}) \\ & s.t. \boldsymbol{\mu} \geq \mathbf{0}, \end{aligned} \tag{11}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g_j : \mathbb{R}^n \rightarrow \mathbb{R}$ are given functions, \mathcal{X} is a subset of \mathbb{R}^n , and

$$q(\boldsymbol{\mu}) = \inf_{\mathbf{x} \in \mathcal{X}} L(\mathbf{x}, \boldsymbol{\mu}) = \inf_{\mathbf{x} \in \mathcal{X}} \{f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x})\}$$

is the dual function.

Algorithms for constrained optimization

- Dual Derivatives and Subgradients - Pros

It is worth reflecting on the potential incentives for solving the dual problem in place of the primal. These are:

- a) The dual is a concave problem (concave cost, convex constraint set). By contrast, the primal need not be convex.
- b) The dual may have smaller dimension and/or simpler constraints than the primal.

Algorithms for constrained optimization

- Dual Derivatives and Subgradients - Pros
- c) If there is no duality gap and the dual is solved exactly to yield a Lagrange multiplier $\boldsymbol{\mu}^*$, all optimal primal solutions can be obtained by minimizing the Lagrangian $L(\mathbf{x}, \boldsymbol{\mu}^*)$ over $\mathbf{x} \in \mathcal{X}$ (however, there may be additional minimizers of $L(\mathbf{x}, \boldsymbol{\mu}^*)$ that are primal-infeasible). Furthermore, if the dual is solved approximately to yield an approximate Lagrange multiplier $\boldsymbol{\mu}$, and \mathbf{x}_μ minimizes $L(\mathbf{x}, \boldsymbol{\mu})$ over $\mathbf{x} \in \mathcal{X}$, then \mathbf{x}_μ also solves

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & s.t. \mathbf{x} \in \mathcal{X}, g_j(\mathbf{x}) \leq g_j(\mathbf{x}_\mu), j = 1, \dots, r. \end{aligned} \tag{10}$$

Thus if the constraint violations $g_j(\mathbf{x}_\mu)$ are not much larger than zero, \mathbf{x}_μ may be an acceptable practical solution.

- d) Even if there is a duality gap, for every $\boldsymbol{\mu} \geq \mathbf{0}$, the dual value $q(\boldsymbol{\mu})$ is a lower bound to the optimal primal value. This lower bound may be useful in the context of discrete optimization and branch and bound procedures.

Algorithms for constrained optimization

- Dual Derivatives and Subgradients - Cons

We should also consider some of the difficulties in solving the dual problem. The most important ones are the following:

- a) To evaluate the dual function at any $\boldsymbol{\mu}$ requires minimization of the Lagrangian $L(\mathbf{x}, \boldsymbol{\mu})$ over $\mathbf{x} \in \mathcal{X}$. In effect, this restricts the utility of dual methods to problems where this minimization can either be done in closed form or else is relatively simple; for example, when there is special structure that allows decomposition, as in the separable problems and the monotropic programming problems.
- b) In many types of problems, the dual function is nondifferentiable, in which algorithms for smooth objective functions do not apply.
- c) Even if we find an optimal dual solution $\boldsymbol{\mu}^*$, it may be difficult to obtain a primal feasible vector \mathbf{x} from the minimization of $L(\mathbf{x}, \boldsymbol{\mu}^*)$ over $\mathbf{x} \in \mathcal{X}$ as required by the primal-dual optimality conditions, since this minimization can also yield primal-infeasible vectors.

Algorithms for constrained optimization

- Dual Derivatives and Subgradients

For a given $\mu \in \mathbb{R}^r$, suppose that \mathbf{x}_μ minimizes the Lagrangian $L(\mathbf{x}, \mu)$,

$$\mathbf{x}_\mu = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} L(\mathbf{x}, \mu) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \{f(\mathbf{x}) + \mu^T \mathbf{g}(\mathbf{x})\}.$$

An important fact is that $\mathbf{g}(\mathbf{x}_\mu)$ is a subgradient of the dual function q at μ :

obtained essentially at no cost

$$q(\bar{\mu}) \leq q(\mu) + (\bar{\mu} - \mu)^T \mathbf{g}(\mathbf{x}_\mu), \quad \forall \bar{\mu} \in \mathbb{R}^r. \quad (10)$$

To see this, we use the definition of q and \mathbf{x}_μ to write for all $\bar{\mu} \in \mathbb{R}^r$,

$$\begin{aligned} q(\bar{\mu}) &= \inf_{\mathbf{x} \in \mathcal{X}} \{f(\mathbf{x}) + \bar{\mu}^T \mathbf{g}(\mathbf{x})\} \leq f(\mathbf{x}_\mu) + \bar{\mu}^T \mathbf{g}(\mathbf{x}_\mu) \\ &= f(\mathbf{x}_\mu) + \mu^T \mathbf{g}(\mathbf{x}_\mu) + (\bar{\mu} - \mu)^T \mathbf{g}(\mathbf{x}_\mu) \\ &= q(\mu) + (\bar{\mu} - \mu)^T \mathbf{g}(\mathbf{x}_\mu). \end{aligned}$$

Note that this calculation is valid for all $\mu \in \mathbb{R}^r$ for which there is a minimizing vector \mathbf{x}_μ , regardless of whether $\mu \geq 0$.

Algorithms for constrained optimization

- Dual Derivatives and Subgradients

Proposition 1. *Let \mathcal{X} be a compact set, and let f and \mathbf{g} be continuous over \mathcal{X} . Assume also that for every $\boldsymbol{\mu} \in \mathbb{R}^r$, $L(\mathbf{x}, \boldsymbol{\mu})$ is minimized over $\mathbf{x} \in \mathcal{X}$ at a unique point \mathbf{x}_μ . Then q is everywhere continuously differentiable and*

$$\nabla q(\boldsymbol{\mu}) = \mathbf{g}(\mathbf{x}_\mu), \quad \forall \boldsymbol{\mu} \in \mathbb{R}^r.$$

Algorithms for constrained optimization

- Lagrangian Algorithms - Equality Constraints

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ s.t. \quad & \mathbf{h}(\mathbf{x}) = \mathbf{0}. \end{aligned}$$

The Lagrangian function is given by

$$l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}). \quad (10)$$

The Lagrangian algorithm for this problem is given by

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha_k (\nabla f(\mathbf{x}^{(k)}) + D\mathbf{h}(\mathbf{x}^{(k)})^T \boldsymbol{\lambda}^{(k)}), \\ \boldsymbol{\lambda}^{(k+1)} &= \boldsymbol{\lambda}^{(k)} + \beta_k \mathbf{h}(\mathbf{x}^{(k+1)}). \end{aligned} \quad (11)$$

The update equation for $\mathbf{x}^{(k)}$ is a gradient algorithm for minimizing the Lagrangian with respect to \mathbf{x} , and the update equation for $\boldsymbol{\lambda}^{(k)}$ is a gradient algorithm for maximizing the Lagrangian with respect to $\boldsymbol{\lambda}$.

Algorithms for constrained optimization

- Lagrangian Algorithms - Equality Constraints

The following lemma establishes that if the algorithm converges, the limit must satisfy the Lagrange condition. More specifically, the lemma states that any fixed point of the algorithm must satisfy the Lagrange condition. A fixed point of an update algorithm is simply a point with the property that when updated using the algorithm, the resulting point is equal to the given point. For the case of the Lagrangian algorithm, which updates both $\mathbf{x}^{(k)}$ and $\boldsymbol{\lambda}^{(k)}$ vectors, a fixed point is a pair of vectors. If the Lagrangian algorithm converges, the limit must be a fixed point.

Lemma 1. *For the Lagrangian algorithm for updating $\mathbf{x}^{(k)}$ and $\boldsymbol{\lambda}^{(k)}$, the pair $(\mathbf{x}^{(k)}, \boldsymbol{\lambda}^{(k)})$ is a fixed point iff it satisfies the Lagrange condition.*

Algorithms for constrained optimization

- Lagrangian Algorithms - Equality Constraints

Below, we use $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ to denote a pair satisfying the Lagrange condition. Assume that $\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \succ \mathbf{0}$. Also assume that \mathbf{x}^* is a regular point. With these assumptions, we are now ready to state and prove that the algorithm is locally convergent. For simplicity, we will take α_k and β_k to be fixed constants (not depending on k), denoted α and β , respectively.

Theorem 1. *For the Lagrangian algorithm for updating $\mathbf{x}^{(k)}$ and $\boldsymbol{\lambda}^{(k)}$, provided that α and β are sufficiently small, there is a neighborhood of $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ such that if the pair $(\mathbf{x}^0, \boldsymbol{\lambda}^0)$ is in this neighborhood, then the the algorithm converges to $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ with at least a linear order of convergence.*

Algorithms for constrained optimization

- Lagrangian Algorithms - Inequality Constraints

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{g}(\mathbf{x}) \leq \mathbf{0}. \end{aligned}$$

The Lagrangian function is given by

$$l(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}). \quad (10)$$

The Lagrangian algorithm for this problem is given by

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha_k (\nabla f(\mathbf{x}^{(k)}) + Dg(\mathbf{x}^{(k)})^T \boldsymbol{\mu}^{(k)}), \\ \boldsymbol{\mu}^{(k+1)} &= [\boldsymbol{\mu}^{(k)} + \beta_k \mathbf{g}(\mathbf{x}^{(k+1)})]_+, \end{aligned}$$

where $[\cdot]_+ = \max\{\cdot, 0\}$. The update equation for $\mathbf{x}^{(k)}$ is a gradient algorithm for minimizing the Lagrangian with respect to its \mathbf{x} argument. The update equation for $\boldsymbol{\mu}^{(k)}$ is a projected gradient algorithm for maximizing the Lagrangian with respect to its $\boldsymbol{\mu}$ argument.

Algorithms for constrained optimization

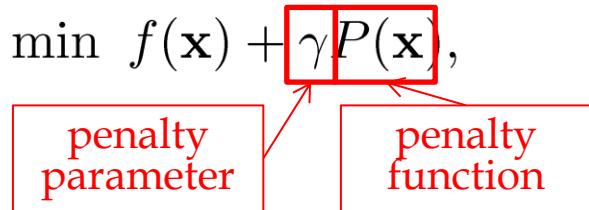
- Penalty Methods

$$\begin{aligned} & \min f(\mathbf{x}) \\ & s.t. \mathbf{x} \in \Omega. \end{aligned} \tag{11}$$

$$\min f(\mathbf{x}) + \gamma P(\mathbf{x}), \tag{12}$$

penalty parameter

penalty function



Definition 1. A function $P : \mathbb{R}^n \rightarrow \mathbb{R}$ is called a *penalty function* for the constrained optimization problem above if it satisfies the following three conditions:

1. P is continuous.
2. $P(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.
3. $P(\mathbf{x}) = 0$ if and only if \mathbf{x} is feasible (i.e., $\mathbf{x} \in \Omega$).

Algorithms for constrained optimization

- Penalty Methods - Choice of penalty function

Example.

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g_i(\mathbf{x}) \leq 0, i = 1, \dots, p. \end{aligned}$$

A possible choice for P is

$$P(\mathbf{x}) = \sum_{i=1}^p g_i^+(\mathbf{x}),$$

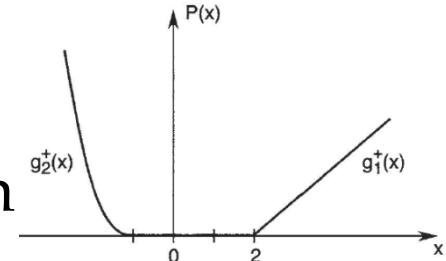
where

$$g_i^+(\mathbf{x}) = \max\{0, g_i(\mathbf{x})\} = \begin{cases} 0, & g_i(\mathbf{x}) \leq 0 \\ g_i(\mathbf{x}), & g_i(\mathbf{x}) > 0. \end{cases}$$

We refer to this penalty function as the absolute value penalty function, because it is equal to $\sum |g_i(\mathbf{x})|$, where the summation is taken over all constraints that are violated at \mathbf{x} .

Algorithms for constrained optimization

- Penalty Methods - Choice of penalty function



Example: Let $g_1, g_2 : \mathbb{R} \rightarrow \mathbb{R}$ be defined by $g_1(x) = x - 2$, $g_2(x) = -(x + 1)^3$. The feasible set defined by $\{x \in \mathbb{R} : g_1(x) \leq 0, g_2(x) \leq 0\}$ is simply the interval $[-1, 2]$. In this example, we have

$$g_1^+(x) = \max\{0, g_1(x)\} = \begin{cases} 0, & x \leq 2 \\ x - 2, & \text{otherwise,} \end{cases}$$

$$g_2^+(x) = \max\{0, g_2(x)\} = \begin{cases} 0, & x \geq -1 \\ -(x + 1)^3, & \text{otherwise,} \end{cases}$$

and

$$P(x) = g_1^+(x) + g_2^+(x) = \begin{cases} x - 2, & x > 2 \\ 0, & -1 \leq x \leq 2 \\ -(x + 1)^3, & x < -1. \end{cases}$$

Algorithms for constrained optimization

- Penalty Methods - Choice of penalty function

The absolute value penalty function may not be differentiable at points \mathbf{x} where $g_i(\mathbf{x}) = 0$. Therefore, in such cases we cannot use techniques for optimization that involve derivatives. A form of the penalty function that is guaranteed to be differentiable is the *Courant-Beltrami penalty function*, given by

$$P(\mathbf{x}) = \sum_{i=1}^p (g_i^+(\mathbf{x}))^2.$$

Algorithms for constrained optimization

- **Penalty Methods - Convergence**

Denote by \mathbf{x}^* a solution (global minimizer) to the problem. Let P be a penalty function for the problem. For each $k = 1, 2, \dots$, let $\gamma_k \in \mathbb{R}$ be a given positive constant. Define an associated function $q(\gamma_k, \cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$q(\gamma_k, \mathbf{x}) = f(\mathbf{x}) + \gamma_k P(\mathbf{x}).$$

For each k , we can write the following associated unconstrained optimization problem:

$$\min q(\gamma_k, \mathbf{x}).$$

Denote by $\mathbf{x}^{(k)}$ a minimizer of $q(\gamma_k, \mathbf{x})$.

Algorithms for constrained optimization

- Penalty Methods - Convergence

Lemma 1. Suppose that $\{\gamma_k\}$ is a nondecreasing sequence; that is, for each k , we have $\gamma_k < \gamma_{k+1}$. Then, for each k we have

1. $q(\gamma_{k+1}, \mathbf{x}^{(k+1)}) \geq q(\gamma_k, \mathbf{x}^{(k)})$.
2. $P(\mathbf{x}^{(k+1)}) \leq P(\mathbf{x}^{(k)})$.
3. $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$.
4. $f(\mathbf{x}^*) \geq q(\gamma_k, \mathbf{x}^{(k)}) \geq f(\mathbf{x}^{(k)})$.

Theorem 2. Suppose that the objective function f is continuous and $\gamma_k \rightarrow \infty$ as $k \rightarrow \infty$. Then, the limit of any convergent subsequence of the sequence $\{\mathbf{x}^{(k)}\}$ is a solution to the constrained optimization problem.

Algorithms for constrained optimization

- Penalty Methods – Exact penalty

We desire an exact solution to the original constrained problem by solving the associated unconstrained problem $\min_{\mathbf{x}} f(\mathbf{x}) + \gamma P(\mathbf{x})$ with a finite $\gamma > 0$. It turns out that indeed this can be accomplished, in which case we say that the penalty function is *exact*. However, it is necessary that exact penalty functions be nondifferentiable.

Example:

$$\begin{aligned} & \min f(x) \\ & s.t. x \in [0, 1], \end{aligned}$$

where $f(x) = 5 - 3x$.

Algorithms for constrained optimization

- Penalty Methods – Exact penalty

Proposition 1. *Consider the problem*

$$\begin{aligned} & \min f(\mathbf{x}) \\ & s.t. \mathbf{x} \in \Omega, \end{aligned}$$

with $\Omega \subset \mathbb{R}^n$ convex. Suppose that the minimizer \mathbf{x}^* lies on the boundary of Ω and there exists a feasible direction \mathbf{d} at \mathbf{x}^* such that $\mathbf{d}^T \nabla f(\mathbf{x}^*) > 0$. If P is an exact penalty function, then P is not differentiable at \mathbf{x}^* .

Proof. We use contraposition. Suppose that P is differentiable at \mathbf{x}^* . Then, $\mathbf{d}^T \nabla P(\mathbf{x}^*) = 0$, because $P(\mathbf{x}) = 0$ for all $\mathbf{x} \in \Omega$. Hence, if we let $g = f + \gamma P$, then $\mathbf{d}^T \nabla g(\mathbf{x}^*) > 0$ for all finite $\gamma > 0$, which implies that $\nabla g(\mathbf{x}^*) \neq 0$. Hence, \mathbf{x}^* is not a local minimizer of g , and thus P is not an exact penalty function. \square