

Chapter 23 Algorithms for Constrained Optimization

1. Projections
2. Projected Gradient Methods with Linear Constraints
3. Lagrangian Algorithms
4. Penalty Methods



Projections

constrained optimization: $\min\{f(\mathbf{x}) \mid \mathbf{x} \in \Omega\}$

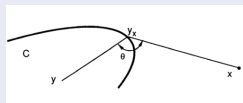
- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is multi-dimensional objective function.
 - $\mathbf{x} \in \Omega \subset \mathbb{R}^n$ is n -dimensional decision variables.
- ★ generic linear search $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$ for unconstrained optimization may not satisfy the constraints.

Definition (projection)

Let $\Omega \subset \mathbb{R}^n$ be nonempty closed convex set.

The projection onto Ω is defined as

$$\Pi_{\Omega}[\mathbf{x}] := \arg \min\{\|\mathbf{z} - \mathbf{x}\| \mid \mathbf{z} \in \Omega\}, \forall \mathbf{x} \in \mathbb{R}^n.$$



properties of projection

- If Ω is nonempty closed convex, then the projection exists and is unique.
- \mathbf{x}^* is the projection of \mathbf{x} onto $\Omega \Leftrightarrow (\mathbf{x} - \mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) \leq 0, \forall \mathbf{z} \in \Omega$. Moreover, “=” holds $\Leftrightarrow \Omega$ is affine.
- projection mapping is nonexpansive, $\|\Pi_{\Omega}[\mathbf{x}] - \Pi_{\Omega}[\mathbf{y}]\| \leq \|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

Projected Methods

projected method for $\min\{f(\mathbf{x}) \mid \mathbf{x} \in \Omega\}$

- (revisit) unconstrained line search: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$.
- projected methods: $\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}]$.
- projected gradient methods: $\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})]$.

projection of peculiar set with explicit solution

- 1 If $\Omega = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{a} \leq \mathbf{x} \leq \mathbf{b}\}$ is a box with $\mathbf{a} \in \mathbb{R}^n$ and $\mathbf{b} \in \mathbb{R}^n$, then $[\Pi_\Omega(\mathbf{x})]_i = \text{median}\{a_i, x_i, b_i\}$ for $i = 1, 2, \dots, n$.
- 2 If $\Omega = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\|_2 \leq a\}$ is a ball with $a > 0$, then $\Pi_\Omega(\mathbf{x}) = \min(1, \frac{a}{\|\mathbf{x}\|_2})\mathbf{x}$.
- 3 If $\Omega = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}\}$ is an affine set with $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, then $\Pi_\Omega(\mathbf{x}) = \mathbf{x} - \mathbf{A}^\dagger(\mathbf{A}\mathbf{x} - \mathbf{b})$, where \mathbf{A}^\dagger is the Moore-Penrose pseudo-inverse of \mathbf{A} . Particularly, if $\Omega = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{a}^\top \mathbf{x} = b, \mathbf{a} \neq \mathbf{0} \in \mathbb{R}^n\}$ is a nonvertical hyplane, then $\Pi_\Omega(\mathbf{x}) = \mathbf{x} - \frac{\mathbf{a}^\top \mathbf{x} - b}{\|\mathbf{a}\|_2^2} \mathbf{a}$.

Projected Gradient Methods (PGM)

Example ($\min \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x}$, s.t. $\|\mathbf{x}\|^2 = 1$, where $\mathbf{Q} = \mathbf{Q}^\top \succ 0$)

- 1 Derive a formula for the update equation for the algorithm.
- 2 Whether the algorithm converges to an optimum or not?
- 3 Show that for $0 < \alpha < \frac{1}{\lambda_{\max}}$ (where λ_{\max} is the largest eigenvalue of \mathbf{Q}), the algorithm converges to an optimum, provided that $\mathbf{x}^{(0)}$ is not orthogonal to the eigenvectors of \mathbf{Q} corresponding to the smallest eigenvalue.

- 1 $\mathbf{x}^{(k+1)} = \beta_k (\mathbf{x}^{(k)} - \alpha \mathbf{Q} \mathbf{x}^{(k)}) = \beta_k (\mathbf{I} - \alpha \mathbf{Q}) \mathbf{x}^{(k)}$, where $\beta_k = \frac{1}{\|(\mathbf{I} - \alpha \mathbf{Q}) \mathbf{x}^{(k)}\|}$.
- 2 If we start with $\mathbf{x}^{(0)}$ being a unit eigenvector of \mathbf{Q} , it is very easy to check that $\mathbf{x}^{(k)} = \mathbf{x}^{(0)}$ for all k .
- 3 Let $\{\mathbf{v}_i\}_{i=1}^n$ be the linear independent eigenvector corresponding to the eigenvalue $\{\lambda_i(\mathbf{Q})\}_{i=1}^n$. We have $\{\mathbf{v}_i\}_{i=1}^n$ is basis of \mathbb{R}^n . Thus,

$$\mathbf{x}^{(k)} = y_1^{(k)} \mathbf{v}_1 + \cdots + y_n^{(k)} \mathbf{v}_n,$$

$$\mathbf{x}^{(k+1)} = y_1^{(k+1)} \mathbf{v}_1 + \cdots + y_n^{(k+1)} \mathbf{v}_n.$$



Projected Gradient Methods (PGM)

It follows from the iterative scheme in item 1 that

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \beta_k(\mathbf{I} - \alpha\mathbf{Q})\mathbf{x}^{(k)} = \beta_k(\mathbf{I} - \alpha\mathbf{Q})\left(y_1^{(k)}\mathbf{v}_1 + \cdots + y_n^{(k)}\mathbf{v}_n\right) \\ &= \beta_k\left(y_1^{(k)}(\mathbf{I} - \alpha\mathbf{Q})\mathbf{v}_1 + \cdots + y_n^{(k)}(\mathbf{I} - \alpha\mathbf{Q})\mathbf{v}_n\right) \\ &= \beta_k\left(y_1^{(k)}(1 - \alpha\lambda_1)\mathbf{v}_1 + \cdots + y_n^{(k)}(1 - \alpha\lambda_n)\mathbf{v}_n\right).\end{aligned}$$

Compared to (2), it yields $y_i^{(k+1)} = \beta_k y_i^{(k)} (1 - \alpha\lambda_i)$.

$$\therefore y_i^{(k)} = (\Pi_{i=0}^{k-1} \beta_k) y_i^{(0)} (1 - \alpha\lambda_i)^k.$$

It follows from (1) that $\mathbf{x}^{(k)} = \sum_{i=1}^n y_i^{(k)} \mathbf{v}_i = y_1^{(k)} \left(\mathbf{v}_1 + \sum_{i=2}^n \frac{y_i^{(k)}}{y_1^{(k)}} \mathbf{v}_i \right)$.

Assume that $y_1^{(0)} \neq 0$, we obtain $\frac{y_i^{(k)}}{y_1^{(k)}} = \frac{y_i^{(0)}(1-\alpha\lambda_i)^k}{y_1^{(0)}(1-\alpha\lambda_1)^k} = \frac{y_i^{(0)}}{y_1^{(0)}} \left(\frac{1-\alpha\lambda_i}{1-\alpha\lambda_1} \right)^k$.

Using the fact that $\frac{1-\alpha\lambda_i}{1-\alpha\lambda_1} < 1$, we deduce that $\frac{y_i^{(k)}}{y_1^{(k)}} \rightarrow 0$,

which implies that $\mathbf{x}^{(k)} \rightarrow \mathbf{v}_1$.



PGM with Linear Constraints

constrained optimization: $\min f(\mathbf{x}), \text{ s.t. } \mathbf{Ax} = \mathbf{b}.$

- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is multi-dimensional objective function.
- $\mathbf{x} \in \mathbb{R}^n$ is n -dimensional decision variables.
- $\mathbf{A} \in \mathbb{R}^{m \times n} (m < n), \text{rank}(\mathbf{A}) = m, \mathbf{b} \in \mathbb{R}^m.$

Lemma

Let $\Omega = \{\mathbf{x} \mid \mathbf{Ax} = \mathbf{b}\}$. Then the projection onto Ω can be defined by the orthogonal projector matrix \mathbf{P} , i.e.,

$$\Pi_{\Omega} = \mathbf{P} = \mathbf{I}_n - \mathbf{A}^{\top}(\mathbf{AA}^{\top})^{-1}\mathbf{A}. \quad (3)$$

key point of proof: By the KKT condition of constrained optimization.



Lemma

For the projection matrix in (3), it holds that $\mathcal{N}(\mathbf{P}) = \mathcal{R}(\mathbf{A}^\top)$, $\mathcal{N}(\mathbf{A}) = \mathcal{R}(\mathbf{P})$.

proof. \Rightarrow) For any $\mathbf{v} \in \mathcal{N}(\mathbf{P})$, we have

$$0 = \mathbf{P}\mathbf{v} = (\mathbf{I}_n - \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A})\mathbf{v} = \mathbf{v} - \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A}\mathbf{v}.$$

$$\therefore \mathbf{v} = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A}\mathbf{v}, \text{ and hence } \mathbf{v} \in \mathcal{R}(\mathbf{A}^\top).$$

\Leftarrow) For any $\mathbf{v} \in \mathcal{R}(\mathbf{A}^\top)$, then there exists $\mathbf{u} \in \mathbb{R}^m$ such that $\mathbf{v} = \mathbf{A}^\top\mathbf{u}$.

$$\begin{aligned} \therefore \mathbf{P}\mathbf{v} &= (\mathbf{I}_n - \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A})\mathbf{A}^\top\mathbf{u} \\ &= \mathbf{A}^\top\mathbf{u} - \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A}\mathbf{A}^\top\mathbf{u} = 0. \end{aligned}$$

Hence, we have proved that $\mathcal{N}(\mathbf{P}) = \mathcal{R}(\mathbf{A}^\top)$.

By arguments similar to the above, we can show $\mathcal{N}(\mathbf{A}) = \mathcal{R}(\mathbf{P})$.



PGM with Linear Constraints

revisit: for the unconstrained optimization $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$, the first-order necessary condition (FONC) for a point \mathbf{x}^* to be a local minimizer is $\nabla f(\mathbf{x}^*) = 0$.

Theorem

For constrained optimization $\min\{f(\mathbf{x}) \mid \mathbf{x} \in \Omega\}$ with $\Omega = \{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{b}\}$, the FONC for a point \mathbf{x}^ to be a local minimizer is $\mathbf{P}\nabla f(\mathbf{x}^*) = 0$.*

proof. $\mathbf{P}\nabla f(\mathbf{x}^*) = 0$,

$$\iff \nabla f(\mathbf{x}^*) \in \mathcal{N}(\mathbf{P}), \quad [\text{by the definition of } \mathcal{N}(\mathbf{P})]$$

$$\iff \nabla f(\mathbf{x}^*) \in \mathcal{R}(\mathbf{A}^\top), \quad [\text{by the above lemma}]$$

$$\iff \exists \boldsymbol{\lambda}^* \in \mathbb{R}^m, \text{ such that } \nabla f(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\lambda}^* = 0,$$

$$\iff \text{by combining } \mathbf{x}^* \in \Omega = \{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{b}\},$$

we can obtain the KKT condition

$$\begin{cases} \nabla f(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\lambda}^* = 0, \\ \mathbf{A}^\top \mathbf{x}^* = \mathbf{b}. \end{cases}$$



PGM with Linear Constraints

iterative scheme of projection gradient method (PGM)

$\min\{f(\mathbf{x}) \mid \mathbf{Ax} = \mathbf{b}\}$, the iterative scheme of PGM is

$$\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})] \stackrel{\text{why?}}{=} \mathbf{x}^{(k)} - \alpha_k \mathbf{P} \nabla f(\mathbf{x}^{(k)}).$$

Theorem

In PGM, if $\mathbf{x}^{(0)}$ is feasible, then each $\mathbf{x}^{(k)}$ is feasible, i.e., $\mathbf{Ax}^{(k)} = \mathbf{b}$, $\forall k \geq 0$.

proof. (by induction).

- The result holds for $k = 0$ by assumption.
- Suppose that $\mathbf{Ax}^{(k)} = \mathbf{b}$. We now show that $\mathbf{Ax}^{(k+1)} = \mathbf{b}$.

$$\because \mathbf{P} \nabla f(\mathbf{x}^{(k)}) \in \mathcal{R}(\mathbf{P}) \stackrel{\text{lemma}}{\implies} \mathbf{P} \nabla f(\mathbf{x}^{(k)}) \in \mathcal{N}(\mathbf{A}).$$

$$\begin{aligned} \therefore \mathbf{Ax}^{(k+1)} &= \mathbf{A}(\mathbf{x}^{(k)} - \alpha_k \mathbf{P} \nabla f(\mathbf{x}^{(k)})) \\ &= \mathbf{Ax}^{(k)} - \alpha_k \mathbf{AP} \nabla f(\mathbf{x}^{(k)}) \\ &= \mathbf{b}. \quad (\text{why? Ans: } \mathbf{AP} \nabla f(\mathbf{x}^{(k)}) = \mathbf{0}) \end{aligned}$$



Properties of PGM with Linear Constraints

Theorem

If $\{\mathbf{x}^{(k)}\}$ is a sequence generated by the projected steepest descent method and if $\mathbf{P}\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$, then $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$.

proof. $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{P}\nabla f(\mathbf{x}^{(k)})$ with $\alpha_k \geq 0$ is the exact step size.

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

Then, $\phi_k(\alpha_k) \leq \phi_k(\alpha)$ for all $\alpha \geq 0$.

By the chain rule,

$$\begin{aligned}\phi'_k(0) &= -\nabla f(\mathbf{x}^{(k)} - 0\mathbf{P}\nabla f(\mathbf{x}^{(k)}))^\top \mathbf{P}\nabla f(\mathbf{x}^{(k)}) \\ &= -\nabla f(\mathbf{x}^{(k)})^\top \mathbf{P}\nabla f(\mathbf{x}^{(k)}) \\ &= -\|\mathbf{P}\nabla f(\mathbf{x}^{(k)})\|^2 < 0. \quad [\because \mathbf{P} = \mathbf{P}^2 = \mathbf{P}^\top \mathbf{P}]\end{aligned}$$

$\therefore \exists \bar{\alpha} > 0$ such that $\phi_k(0) > \phi_k(\alpha)$, $\forall \alpha \in (0, \bar{\alpha}]$.

$\therefore f(\mathbf{x}^{(k+1)}) = \phi_k(\alpha_k) \leq \phi_k(\bar{\alpha}) < \phi_k(0) = f(\mathbf{x}^{(k)})$.



PGM for Linear Program

$$\text{(LP)} \quad \min \mathbf{c}^\top \mathbf{x}, \quad \text{s.t. } \mathbf{Ax} \geq \mathbf{b}, \text{ where } \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m \text{ and } \mathbf{c} \in \mathbb{R}^n.$$

Lemma

Let the solution set of (LP) be nonempty, then for any $\mathbf{x}_0 \in \mathbb{R}^n$, there exists $\varepsilon_0 > 0$, such that for all $\varepsilon \in (0, \varepsilon_0]$,

$$\text{(LP)} \iff \min \frac{\varepsilon}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \mathbf{c}^\top \mathbf{x}, \quad \text{s.t. } \mathbf{Ax} \geq \mathbf{b}.$$

Proof: By using the KKT condition.



PGM for Linear Program

discussion

$$\min \mathbf{c}^\top \mathbf{x}, \quad \text{s.t. } \mathbf{Ax} \geq \mathbf{b}.$$



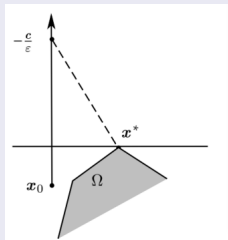
$$\min \frac{\varepsilon}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \mathbf{c}^\top \mathbf{x}, \quad \text{s.t. } \mathbf{Ax} \geq \mathbf{b}.$$



$$\min \frac{\varepsilon}{2} \left\| \mathbf{x} - \left(\mathbf{x}_0 - \frac{\mathbf{c}}{\varepsilon} \right) \right\|^2, \quad \text{s.t. } \mathbf{Ax} \geq \mathbf{b}.$$



Projection of $\left(\mathbf{x}_0 - \frac{\mathbf{c}}{\varepsilon} \right)$ onto the feasible set $\Omega = \{\mathbf{x} \mid \mathbf{Ax} \geq \mathbf{b}\}$.



★ For any initial point \mathbf{x}_0 , as long as the step size is large enough for projected gradient algorithm, the optimum is obtained in one step (Mangasarian, 1979).



Lagrangian Algorithms for Equality Constraints

motivation for Lagrangian algorithms

By using gradient algorithms to update simultaneously the decision variable and Lagrange multiplier vector.

equality constrained optimization

$$\min f(\mathbf{x}), \text{ s.t. } \mathbf{h}(\mathbf{x}) = 0, \text{ where } \mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m. \quad (4)$$

Lagrangian function of (4) is $l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x})$.

Iterative scheme of Lagrangian method for (4) is

$$\begin{cases} \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla_{\mathbf{x}} l(\mathbf{x}^{(k)}, \boldsymbol{\lambda}^{(k)}) = \mathbf{x}^{(k)} - \alpha_k (\nabla f(\mathbf{x}^{(k)}) + \mathbf{D}\mathbf{h}(\mathbf{x}^{(k)})^\top \boldsymbol{\lambda}^{(k)}), \\ \boldsymbol{\lambda}^{(k+1)} = \boldsymbol{\lambda}^{(k)} + \beta_k \nabla_{\boldsymbol{\lambda}} l(\mathbf{x}^{(k)}, \boldsymbol{\lambda}^{(k)}) = \boldsymbol{\lambda}^{(k)} + \beta_k \mathbf{h}(\mathbf{x}^{(k)}). \end{cases} \quad (5)$$

★ $\mathbf{x}^{(k)}$ is a **gradient method** for **minimizing** Lagrangian w.r.t \mathbf{x} variable;
 $\boldsymbol{\lambda}^{(k)}$ is a **gradient method** for **maximizing** Lagrangian w.r.t $\boldsymbol{\lambda}$ variable.



Lagrangian Algorithms for Equality Constraints

Lemma

When Lagrangian method (5) updates $\mathbf{x}^{(k)}$ and $\boldsymbol{\lambda}^{(k)}$, the pair $(\mathbf{x}^, \boldsymbol{\lambda}^*)$ is a fixed point if and only if it satisfies the Lagrange condition.*

Theorem

When Lagrangian method (5) updates $\mathbf{x}^{(k)}$ and $\boldsymbol{\lambda}^{(k)}$, provided that α_k and β_k are sufficiently small, there is a neighborhood of $(\mathbf{x}^, \boldsymbol{\lambda}^*)$ such that: if $(\mathbf{x}^{(0)}, \boldsymbol{\lambda}^{(0)})$ is in this neighborhood, then the method converges to $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ with at least a linear order of convergence.*

(Proof: By using the contraction mapping fixed point theorem.)



Lagrangian Algorithms for Inequality Constraints

inequality constrained optimization

$$\min f(\mathbf{x}), \quad \text{s.t. } \mathbf{g}(\mathbf{x}) \leq 0, \quad \text{where } \mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p. \quad (6)$$

Lagrangian function of (6): $l(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x})$, $\boldsymbol{\mu} \geq 0$.

Lagrangian method for (6):

$$\begin{cases} \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k (\nabla f(\mathbf{x}^{(k)}) + \mathbf{Dg}(\mathbf{x}^{(k)})^\top \boldsymbol{\mu}^{(k)}), \\ \boldsymbol{\mu}^{(k+1)} = [\boldsymbol{\mu}^{(k)} + \beta_k \mathbf{g}(\mathbf{x}^{(k)})]_+, \quad \text{where } [\cdot]_+ = \max\{\cdot, 0\}. \end{cases} \quad (7)$$

- ★ $\mathbf{x}^{(k)}$ is a **gradient method** for **minimizing** Lagrangian w.r.t its \mathbf{x} variable;
 $\boldsymbol{\mu}^{(k)}$ is a **projected gradient method** for **maximizing** Lagrangian w.r.t $\boldsymbol{\mu}$.



Lagrangian Algorithms for Inequality Constraints

Lemma

When Lagrangian method (7) updates $\mathbf{x}^{(k)}$ and $\boldsymbol{\mu}^{(k)}$, the pair $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a fixed point if and only if it satisfies the KKT condition.

Theorem

When Lagrangian method (7) updates $\mathbf{x}^{(k)}$ and $\boldsymbol{\mu}^{(k)}$, provided that α_k and β_k are sufficiently small, there is a neighborhood of $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ such that if the pair $(\mathbf{x}^{(0)}, \boldsymbol{\mu}^{(0)})$ is in this neighborhood, then

- 1 the nonactive multipliers reduce to zero in finite time and remain at zero thereafter;
- 2 the algorithm converges to $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ with at least a linear order of convergence.

proof: by using contraction mapping fixed point theorem and nonexpansive mapping theorem.



Penalty Method

motivation for penalty method

The constrained optimization is approximated and solved by a sequence of unconstrained optimization.

Definition

indicator function Let $\Omega \subset \mathbb{R}^n$ be nonempty. The indicator function is defined by

$$\iota_{\Omega}(\mathbf{x}) := \begin{cases} 0, & \text{if } \mathbf{x} \in \Omega. \\ +\infty, & \text{otherwise.} \end{cases}$$

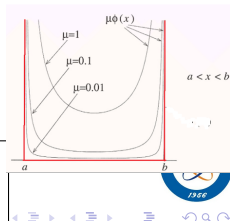
constrained optimization: $\min f(\mathbf{x}), \text{ s.t. } \mathbf{x} \in \Omega.$

\Updownarrow by definition of ι_{Ω}

unconstrained optimization: $\min f(\mathbf{x}) + \iota_{\Omega}(\mathbf{x}).$

\Downarrow approximated by following problem

unconstrained optimization: $\min f(\mathbf{x}) + \gamma \mathbf{P}(\mathbf{x}),$
where $\gamma > 0$ is a constant; $\mathbf{P} : \mathbb{R}^n \rightarrow \mathbb{R}$ is a given function.



Penalty Method

$$(I) \min_{\mathbf{x} \in \Omega} f(\mathbf{x}),$$

\approx

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

where $\gamma > 0$ is a constant;
 $\mathbf{P} : \mathbb{R}^n \rightarrow \mathbb{R}$ is a given function.

★ γ is called the penalty parameter, \mathbf{P} is called the penalty function.

Definition (penalty function)

A function $\mathbf{P} : \mathbb{R}^n \mapsto \mathbb{R}$ is called a penalty function for the constrained optimization if it satisfies the following three conditions:

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

★ The accuracy of (II) for approximating (I) depends on γ and \mathbf{P} .

★ The larger the γ , the closer the approximation. Ideally, as $\gamma \rightarrow +\infty$, the (II) should yield the true solution to (I).



Penalty Methods

Given $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^p$, $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

constrained optimization: $\min f(\mathbf{x})$, s.t. $\mathbf{h}(\mathbf{x}) = 0$, $\mathbf{g}(\mathbf{x}) \leq 0$.

- exact penalty function:

$$\min f(\mathbf{x}) + \gamma \left(\sum_{i=1}^p |\mathbf{h}_i(\mathbf{x})| + \sum_{i=1}^m |\mathbf{g}_i^+(\mathbf{x})| \right),$$

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

- quadratic penalty function: $\min f(\mathbf{x}) + \gamma \left(\sum_{i=1}^p |\mathbf{h}_i(\mathbf{x})|^2 + \sum_{i=1}^m |\mathbf{g}_i^+(\mathbf{x})|^2 \right)$.
- Generalized penalty function:
 $\min f(\mathbf{x}) + \mathbf{P}(\mathbf{x}) = f(\mathbf{x}) + \gamma \left(\|\mathbf{h}(\mathbf{x})\|_b^a + \|[\mathbf{g}(\mathbf{x})]^+\|_b^a \right), \quad a > 0, b > 0.$

★ exact penalty function, nondifferentiable; quadratic penalty, generally differentiable.



constraint $\Omega = \{x \in \mathbb{R} \mid g_1(x) = x - 2 \leq 0, g_2(x) = -(x + 1)^3 \leq 0\}$

Ans: Indeed, $\Omega = [-1, 2]$.

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

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

Thus, the exact penalty function is

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



Penalty Methods for Equality Constraints

Example (solve by quadratic penalty function method)

$$\min f(\mathbf{x}) = -x_1x_2, \text{ s.t. } g(\mathbf{x}) = x_1 + 2x_2 - 4 = 0.$$

Ans: quadratic penalty function is

$$\min_{\mathbf{x}} q(\mathbf{x}) = -x_1x_2 + \frac{1}{2}\gamma (x_1 + 2x_2 - 4)^2. \quad (8)$$

By solving unconstrained optimization (8) with parameter $\gamma > 0$, we have

$$\begin{cases} -x_2 + \gamma (x_1 + 2x_2 - 4) = 0 \\ -x_1 + 2\gamma (x_1 + 2x_2 - 4) = 0 \end{cases} \Rightarrow \begin{cases} x_1(\gamma) = \frac{8\gamma}{4\gamma-1} \\ x_2(\gamma) = \frac{4\gamma}{4\gamma-1} \end{cases} \xrightarrow{\gamma \rightarrow +\infty} \begin{cases} x_1 = 2 \\ x_2 = 1 \end{cases}$$

By comparing the optimality condition of (8) and KKT of original constrained optimization, it yields a dual solution

$$\lambda(\gamma) \stackrel{\text{why?}}{=} -\gamma g(\mathbf{x}(\gamma)) = \frac{-4\gamma}{4\gamma-1} \xrightarrow{\gamma \rightarrow +\infty} \lambda = -1.$$



Penalty Methods for Equality Constraints

$$\min \mathbf{x}^\top \mathbf{Q} \mathbf{x}, \quad \text{s.t. } \|\mathbf{x}\|^2 = 1, \quad \text{where } \mathbf{Q} = \mathbf{Q}^\top.$$

- ① write the quadratic penalty function with penalty parameter γ .
- ② show that \mathbf{x}_γ is an eigenvector of \mathbf{Q} for any γ .
- ③ show that $\|\mathbf{x}_\gamma\|^2 - 1 = O(1/\gamma)$, as $\gamma \rightarrow +\infty$.

Ans: quadratic penalty function: $\min q(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \gamma (\|\mathbf{x}\|^2 - 1)^2$.

As \mathbf{x}_γ is a minimizer of $q(\mathbf{x})$, it follows from optimality condition for unconstrained optimization that

$$\nabla q(\mathbf{x}_\gamma) = 0 \implies 2\mathbf{Q}\mathbf{x}_\gamma + 4\gamma (\|\mathbf{x}_\gamma\|^2 - 1) \mathbf{x}_\gamma = 0,$$

$$\text{i.e., } \underbrace{\mathbf{Q}\mathbf{x}_\gamma = 2\gamma (1 - \|\mathbf{x}_\gamma\|^2) \mathbf{x}_\gamma}_{=:\lambda_\gamma} =: \lambda_\gamma \mathbf{x}_\gamma, \quad \therefore \mathbf{x}_\gamma \text{ is an eigenvector of } \mathbf{Q}.$$

$$\lambda_\gamma = 2\gamma (1 - \|\mathbf{x}_\gamma\|^2) \leq \lambda_{\max},$$

$$\therefore \|\mathbf{x}_\gamma\|^2 - 1 = -\frac{\lambda_{\max}}{2\gamma} = O\left(\frac{1}{\gamma}\right) \text{ as } \gamma \rightarrow \infty.$$



Penalty Methods for Inequality Constraints

solve constrained optimization by quadratic penalty function method

$$\begin{array}{ll}\min & f(\mathbf{x}) = x_1 + x_2, \\ \text{s.t.} & g_1(\mathbf{x}) = -x_1^2 + 2x_2 \geq 0, \quad g_2(\mathbf{x}) = x_1 \geq 0.\end{array}$$

Ans: quadratic penalty function is

$$\min q(\mathbf{x}) = f(\mathbf{x}) + \gamma [\min(0, -x_1^2 + x_2)]^2 + \gamma [\min(0, x_1)]^2.$$

By solving the above unconstrained optimization problem with parameter γ ,

$$\begin{cases} 1 + 2\gamma[\min(0, -x_1^2 + x_2)](-2x_1) + 2\gamma[\min(0, x_1)] = 0, \\ 1 + 2\gamma[\min(0, -x_1^2 + x_2)] = 0, \end{cases}$$

$$\Rightarrow \begin{cases} 1 + 4\gamma x_1 (x_1^2 - x_2) + 2\gamma x_1 = 0 \\ 1 - 2\gamma (x_1^2 - x_2) = 0 \end{cases} \Rightarrow \begin{cases} x_1 = -\frac{1}{2+2\gamma} \\ x_2 = \frac{1}{(2+2\gamma)^2} - \frac{1}{2\gamma} \end{cases}$$

$$\xrightarrow{\gamma \rightarrow +\infty} \begin{cases} x_1 = 0. \\ x_2 = 0. \end{cases} \quad (\text{check the solution by KKT})$$



Penalty Methods

constrained optimization with only inequality constraints

$\min f(\mathbf{x})$, s.t. $\mathbf{g}(\mathbf{x}) \geq 0$, where $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

① logarithmic penalty function: $\min q(\mathbf{x}) = f(\mathbf{x}) - \frac{1}{\gamma} \left(\sum_{i=1}^m \log(\mathbf{g}_i(\mathbf{x})) \right)$.

② reciprocal penalty function: $\min q(\mathbf{x}) = f(\mathbf{x}) - \frac{1}{\gamma} \left(\sum_{i=1}^m \frac{1}{\mathbf{g}_i(\mathbf{x})} \right)$.

solve constrained optimization by logarithmic penalty function method

$\min f(\mathbf{x}) = x_1 - 2x_2$, s.t. $1 + x_1 - x_2^2 \geq 0$, $x_2 \geq 0$.

Ans: logarithmic penalty function is

$$\min q(\mathbf{x}) = f(\mathbf{x}) - \frac{1}{\gamma} [\log(1 + x_1 - x_2^2) + \log(x_2)].$$

By solving the above unconstrained optimization

$$\begin{cases} 1 - \frac{1}{\gamma(1+x_1-x_2^2)} = 0 \\ -2 + \frac{2x_2}{\gamma(1+x_1-x_2^2)} - \frac{1}{\gamma x_2} = 0 \end{cases} \Rightarrow \begin{cases} x_1(\gamma) = \frac{\sqrt{\gamma^2+2\gamma+3}-\gamma}{2\gamma} \\ x_2(\gamma) = \frac{\gamma+\sqrt{\gamma^2+2\gamma}}{2\gamma} \end{cases} \xrightarrow{\gamma \uparrow +\infty} \begin{cases} x_1 = 0 \\ x_2 = 1 \end{cases}$$



Penalty Methods

Definition (associated function)

Let $\{\gamma_k\}_{k=0}^{\infty}$ be a positive sequence. The associated function $q(\gamma_k, \cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$, is defined by

$$q(\gamma_k, \mathbf{x}) = f(\mathbf{x}) + \gamma_k P(\mathbf{x}). \quad (9)$$

★ γ_k are dynamic penalty factors. Let $\mathbf{x}^{(k)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} q(\gamma_k, \mathbf{x})$.

Lemma

Suppose the sequence $\{\gamma_k\}$ is positive and nondecreasing. Then,

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

- ① By the definition of q in (9), the monotonicity of $\{\gamma_k\}$ and $\mathbf{x}^{(k)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} q(\gamma_k, \mathbf{x})$, we have

$$\begin{aligned} q(\gamma_{k+1}, \mathbf{x}^{(k+1)}) &= f(\mathbf{x}^{(k+1)}) + \gamma_{k+1} P(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k+1)}) \\ &\geq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k)}) = q(\gamma_k, \mathbf{x}^{(k)}). \end{aligned}$$



Penalty Methods

- ② $\because \mathbf{x}^{(k)} = \arg \min q(\gamma_k, \mathbf{x})$ and $\mathbf{x}^{(k+1)} = \arg \min q(\gamma_{k+1}, \mathbf{x})$. Thus,
- $$\begin{cases} q(\gamma_k, \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}) \leq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k+1)}) \\ q(\gamma_{k+1}, \mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k+1)}) + \gamma_{k+1} P(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) + \gamma_{k+1} P(\mathbf{x}^{(k)}) \end{cases}$$
- By adding the above inequalities, it yields

$$\gamma_k P(\mathbf{x}^{(k)}) + \gamma_{k+1} P(\mathbf{x}^{(k+1)}) \leq \gamma_{k+1} P(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k+1)}).$$

By rearranging terms, we get $P(\mathbf{x}^{(k+1)}) \leq P(\mathbf{x}^{(k)})$.

- ③ $\because \mathbf{x}^{(k)} = \arg \min q(\gamma_k, \mathbf{x})$
- $$\begin{aligned} \therefore q(\gamma_k, \mathbf{x}^{(k)}) &= f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}) \leq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k+1)}). \\ \therefore f(\mathbf{x}^{(k+1)}) &\geq f(\mathbf{x}^{(k)}) + \gamma_k (P(\mathbf{x}^{(k)}) - P(\mathbf{x}^{(k+1)})). \end{aligned}$$
- By item 2, we get $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$.

- ④ $\because \mathbf{x}^{(k)} = \arg \min q(\gamma_k, \mathbf{x})$
- $$\begin{aligned} \therefore f(\mathbf{x}^*) + \gamma_k P(\mathbf{x}^*) &\geq q(\gamma_k, \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}). \\ \therefore \mathbf{x}^* &\text{ is a minimizer of constrained optimization problem,} \\ \therefore P(\mathbf{x}^*) &= 0 \text{ and } f(\mathbf{x}^*) \geq f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}). \\ \therefore P(\mathbf{x}^{(k)}) &\geq 0 \text{ and } \gamma_k \geq 0, \\ \therefore f(\mathbf{x}^*) &\geq q(\gamma_k, \mathbf{x}^{(k)}) \geq f(\mathbf{x}^{(k)}). \end{aligned}$$



Theorem

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

proof: suppose $\{\mathbf{x}^{(m_k)}\}$ is a subsequence of $\{\mathbf{x}^{(k)}\}$ converging to $\hat{\mathbf{x}}$.

By above Lemma, $\{q(\gamma_k, \mathbf{x}^{(k)})\}$ is nondecreasing and bounded above by $f(\mathbf{x}^*)$.

$\therefore \{q(\gamma_k, \mathbf{x}^{(k)})\}$ has a limit, i.e., $q^* = \lim_{k \rightarrow \infty} q(\gamma_k, \mathbf{x}^{(k)})$ such that $q^* \leq f(\mathbf{x}^*)$.

$\therefore f$ is continuous and $f(\mathbf{x}^{(m_k)}) \leq f(\mathbf{x}^*)$ by the above Lemma,

$$\lim_{k \rightarrow \infty} f(\mathbf{x}^{(m_k)}) = f(\lim_{k \rightarrow \infty} \mathbf{x}^{(m_k)}) = f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*).$$

\therefore both $\{f(\mathbf{x}^{(m_k)})\}$ and $\{q(\gamma_{m_k}, \mathbf{x}^{(m_k)})\}$ converge,



Penalty Methods

$\therefore \{\gamma_{m_k} \mathbf{P}(\mathbf{x}^{(m_k)})\} = \{q(\gamma_{m_k}, \mathbf{x}^{(m_k)}) - f(\mathbf{x}^{(m_k)})\}$ also converges, with

$$\lim_{k \rightarrow \infty} \gamma_{m_k} \mathbf{P}(\mathbf{x}^{(m_k)}) = q^* - f(\hat{\mathbf{x}}).$$

By the above Lemma, $\{\mathbf{P}(\mathbf{x}^{(k)})\}$ is nonincreasing and bounded from below by 0.

$\therefore \{\mathbf{P}(\mathbf{x}^{(k)})\}$ converges, and hence so does $\{\mathbf{P}(\mathbf{x}^{(m_k)})\}$.

$\therefore \gamma_{m_k} \rightarrow \infty, \therefore \lim_{k \rightarrow \infty} \mathbf{P}(\mathbf{x}^{(m_k)}) = 0$. By continuity of $\mathbf{P}(\mathbf{x})$, we have

$$0 = \lim_{k \rightarrow \infty} \mathbf{P}(\mathbf{x}^{(m_k)}) = \mathbf{P}(\lim_{k \rightarrow \infty} \mathbf{x}^{(m_k)}) = \mathbf{P}(\hat{\mathbf{x}}),$$

$\therefore \hat{\mathbf{x}}$ is a feasible point.

$\therefore f(\mathbf{x}^*) \geq f(\hat{\mathbf{x}})$ from above,

$\therefore \hat{\mathbf{x}}$ must be a solution to the constrained optimization problem.

Conclusion: If $\gamma \rightarrow \infty$, the limit of any convergent subsequence is a minimizer to the original constrained optimization problem.



Penalty Methods

★ (exact penalty functions): To find a minimizer to the original constrained optimization by a single (not a sequence) unconstrained optimization approximating the original problem. It requires an exact solution to the original constrained optimization by solving the unconstrained problem with a finite $\gamma > 0$.

Drawback of exact penalty functions: nondifferentiable.

$$\min f(x) = 5 - 3x, \text{ s.t. } 0 \leq x \leq 1$$

The minimizer is $x^* = 1$.

Suppose that the penalty function P is differentiable at $x^* = 1$.

$\because P(x) = 0$ for all $x \in [0, 1]$ (why?), $P'(x^*) = 0$.

If $g := f + \gamma P$, then

$$g'(x^*) = f'(x^*) + \gamma P'(x^*) \neq 0, \forall \gamma > 0.$$

$\therefore x^* = 1$ does not satisfy the optimality condition to be a local minimizer of g .

$\therefore P$ is not an exact penalty function.

Theorem ($\min \{f(\mathbf{x}) \mid \mathbf{x} \in \Omega\}$ with $\Omega \subset \mathbb{R}^n$ as a convex set)

Suppose that the minimizer \mathbf{x}^ lies on $\text{bd}(\Omega)$ and there exists a feasible direction \mathbf{d} at \mathbf{x}^* such that $\mathbf{d}^\top \nabla f(\mathbf{x}^*) > 0$. If \mathbf{P} is an exact penalty function, then \mathbf{P} is nondifferentiable at \mathbf{x}^* .*

proof. (contradiction) Suppose that \mathbf{P} is differentiable at \mathbf{x}^* .

$$\because P(\mathbf{x}) = 0 \text{ for all } \mathbf{x} \in \Omega,$$

$$\therefore \nabla P(\mathbf{x}^*) = 0.$$

Then, $\mathbf{d}^\top \nabla P(\mathbf{x}^*) = 0$ for all \mathbf{d} .

If $\mathbf{g} := f + \gamma \mathbf{P}$, then $\mathbf{d}^\top \nabla \mathbf{g}(\mathbf{x}^*) = \mathbf{d}^\top \nabla f(\mathbf{x}^*) > 0$ for all $\gamma > 0$, which implies that $\nabla \mathbf{g}(\mathbf{x}^*) \neq \mathbf{0}$.

$\therefore \mathbf{x}^*$ is not a local minimizer of \mathbf{g} ,

$\therefore \mathbf{P}$ is not an exact penalty function.



Augmented Lagrangian Methods

motivation of augmented Lagrangian method (ALM)

By combining Lagrange function with penalty function to avoid large penalty parameters, i.e., use a constant penalty parameter.

Given $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^p$, $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

constrained optimization: $\min f(\mathbf{x})$, s.t. $\mathbf{h}(\mathbf{x}) = 0$, $\mathbf{g}(\mathbf{x}) \leq 0$.

Augmented Lagrangian function $L : \mathbb{R}^n \times \mathbb{R}^p \times \mathbb{R}_+^m \rightarrow \mathbb{R}$,

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x}) + \boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x}) + \gamma(\|\mathbf{h}(\mathbf{x})\|^2 + \|\mathbf{g}^+(\mathbf{x})\|^2)$$

$$\begin{aligned} &= f(\mathbf{x}) + \sum_{i=1}^p \lambda_i h_i(\mathbf{x}) + \sum_{i=1}^m \mu_i g_i(\mathbf{x}) \\ &\quad + \gamma \left(\sum_{i=1}^p |h_i(\mathbf{x})|^2 + \sum_{i=1}^m |g_i^+(\mathbf{x})|^2 \right). \end{aligned}$$



iterative scheme of ALM

$$\begin{cases} \mathbf{x}^{k+1} = \arg \min L(\mathbf{x}, \boldsymbol{\lambda}^k, \boldsymbol{\mu}^k), \\ \boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k + \gamma \sum_{i=1}^p h_i(\mathbf{x}^{k+1}), \\ \boldsymbol{\mu}^{k+1} = \left(\boldsymbol{\mu}^k + \gamma \sum_{i=1}^m |g_i^+(\mathbf{x})|^2 \right)^+. \end{cases}$$

