# MATH 4211/6211 – Optimization

# Algorithms for Constrained Optimization

Xiaojing Ye

Department of Mathematics & Statistics

Georgia State University

We know that the gradient method proceeds as

$$x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$$

where $d^{(k)}$ is a descent direction (often chosen as a function of $g^{(k)}$).

However, $x^{(k+1)}$ is not necessarily in the feasible set $\Omega$.

Hence the *projected gradient* (PG) method proceeds as

$$x^{(k+1)} = \Pi(x^{(k)} + \alpha_k d^{(k)})$$

in order that $x^{(k)} \in \Omega$ for all $k$. Here $\Pi(x)$ is the *projection of x onto $\Omega$*.

**Definition.** The projection $\Pi$ onto $\Omega$ is defined by

$$\Pi(z) = \arg\min_{x \in \Omega} \|x - z\|$$

Namely, $\Pi(x)$ is the "closest point" in $\Omega$ to $x$.

Note that $\Pi(x)$ is itself an optimization problem, which may not have closed-form or be easy to solve in most cases.

**Example.** Find the projection operators $\Pi(x)$ for the following sets $\Omega$:

1. $\Omega = \{x \in \mathbb{R}^n : \|x\|_\infty \leq 1\}$

2. $\Omega = \{x \in \mathbb{R}^n : a_i \leq x_i \leq b_i, \ \forall\, i\}$

3. $\Omega = \{x \in \mathbb{R}^n : \|x\| \leq 1\}$

4. $\Omega = \{x \in \mathbb{R}^n : \|x\| = 1\}$

5. $\Omega = \{x \in \mathbb{R}^n : \|x\|_1 \leq 1\}$

6. $\Omega = \{x \in \mathbb{R}^n : \boldsymbol{A}\boldsymbol{x} = \boldsymbol{0}\}$ where $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ with $m \leq n$ is full rank.

**Example.** Consider the constrained optimization problem:

$$\text{minimize} \quad \frac{1}{2}x^\top Q x$$
$$\text{subject to} \quad \|x\|^2 = 1$$

where $Q \succ 0$. Apply the PG method with a fixed step size $\alpha > 0$ to this problem. Specifically:

- Write down the explicit formula of $x^{(k+1)}$ in terms of $x^{(k)}$ (assume never projecting $0$).

- Is it possible to ensure convergence when $\alpha$ is sufficiently small?

- Show that if $\alpha \in (0, \frac{1}{\lambda_{\max}})$ and $x^{(0)}$ is not orthogonal to the smallest eigenvector corresponding to $\lambda_{\min}$, then $x^{(k)}$ converges. Here $\lambda_{\max}$ ($\lambda_{\min}$) is the largest (smallest) eigenvalue of $Q$.

**Solution.** We can see that the solution should be a unit eigenvector corresponding to $\lambda_{\min}$.

Recall that $\Pi(x) = \frac{x}{\|x\|}$ for all $x \neq 0$.

We also know $\nabla f(x) = Qx$, and $x^{(k)} - \alpha \nabla f(x^{(k)}) = (I - \alpha Q)x^{(k)}$.

Therefore, PG with step size $\alpha$ is given by

$$x^{(k+1)} = \beta_k(I - \alpha Q)x^{(k)}, \quad \text{where } \beta_k = \frac{1}{\|(I - \alpha Q)x^{(k)}\|}$$

Note that, if $x^{(0)}$ is an eigenvector of $Q$ corresponding to eigenvalue $\lambda$, then

$$x^{(1)} = \beta_0(I - \alpha Q)x^{(0)} = \beta_0(1 - \alpha\lambda)x^{(0)} = x^{(0)}$$

and hence $x^{(k)} = x^{(0)}$ for all $k$.

**Solution (cont.)** Denote $\lambda_1 \leq \cdots \leq \lambda_n$ the eigenvalues of $\boldsymbol{Q}$, and $\boldsymbol{v}_1, \ldots, \boldsymbol{v}_n$ the corresponding eigenvectors.

Now assume that

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

Then we have

$$\boldsymbol{x}^{(k+1)} = \Pi((\boldsymbol{I} - \alpha \boldsymbol{Q})\boldsymbol{x}^{(k)}) = \beta_k y_1^{(k)}(1 - \alpha\lambda_1)\boldsymbol{v}_1 + \cdots + \beta_k y_n^{(k)}(1 - \alpha\lambda_n)\boldsymbol{v}_n$$

Denote $\beta^{(k)} = \Pi_{j=0}^{k-1} \beta_j$, then

$$y_i^{(k)} = \beta_{k-1} y_i^{(k-1)}(1 - \alpha\lambda_i) = \cdots = \beta^{(k)} y_i^{(0)}(1 - \alpha\lambda_i)^k$$

**Solution (cont.)** Therefore, we have

$$\boldsymbol{x}^{(k)} = \sum_{i=1}^{n} y_i^{(k)} \boldsymbol{v}_i = y_1^{(k)} \left( \boldsymbol{v}_1 + \sum_{i=2}^{n} \frac{y_i^{(k)}}{y_1^{(k)}} \boldsymbol{v}_i \right)$$

Furthermore,

$$\frac{y_i^{(k)}}{y_1^{(k)}} = \frac{\beta^{(k)} y_i^{(0)} (1 - \alpha \lambda_i)^k}{\beta^{(k)} y_1^{(0)} (1 - \alpha \lambda_1)^k} = \frac{y_i^{(0)}}{y_1^{(0)}} \left( \frac{1 - \alpha \lambda_i}{1 - \alpha \lambda_1} \right)^k$$

Note that $y_1^{(0)} \neq 0$ (since $\boldsymbol{x}^{(0)}$ is not orthogonal to the eigenvector corresponding to $\lambda_1$). As $0 < \alpha < \frac{1}{\lambda_n}$, we have

$$0 < \frac{1 - \alpha \lambda_i}{1 - \alpha \lambda_1} < 1 \quad \Rightarrow \quad \left( \frac{1 - \alpha \lambda_i}{1 - \alpha \lambda_1} \right)^k \to 0 \text{ as } k \to \infty$$

for all $\lambda_i > \lambda_1$. Hence $\boldsymbol{x}^{(k)} \to \boldsymbol{v}_1$.

Projected gradient (PG) method for optimization with linear constraint:

$$\text{minimize} \quad f(\boldsymbol{x})$$
$$\text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b}$$

Then PG is given by

$$\boldsymbol{x}^{(k+1)} = \Pi(\boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)}))$$

where $\Pi$ is the projection onto $\Omega := \{\boldsymbol{x} \in \mathbb{R}^n : \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b}\}$.

We first consider the *orthogonal projection* onto the hyperplane $\Psi = \{x \in \mathbb{R}^n : Ax = 0\}$:

For any $v \in \mathbb{R}^n$, the projection onto $\Psi$ is the solution to

$$\text{minimize} \quad \frac{1}{2}\|x - v\|^2$$
$$\text{subject to} \quad Ax = 0$$

Let $P : \mathbb{R}^n \to \mathbb{R}^n$ denote this projector, i.e., $Pv$ is the point on $\Psi$ closest to $v$.

The Lagrange function is

$$l(x, \lambda) = \frac{1}{2}\|x - v\|^2 + \lambda^\top A x$$

Hence the Lagrange (KKT) condition is

$$(x - v) + A^\top \lambda = 0$$
$$Ax = 0$$

Left-multiplying the first equation by $A$ and using $Ax = 0$, we obtain

$$\lambda = (AA^\top)^{-1} A v$$
$$x = (I - A^\top (AA^\top)^{-1} A) v$$

Denote the projector onto $\Psi$ by

$$P = I - A^\top (AA^\top)^{-1} A$$

Thus, the projection of $v$ onto $\Psi$ is $Pv$.

**Proposition.** The projector $P$ has the following properties:

1. $P = P^\top$

2. $P^2 = P$.

3. $Pv = 0$ iff $\exists \lambda \in \mathbb{R}^m$ s.t. $v = A^\top \lambda$. Namely $\mathcal{N}(P) = \mathcal{R}(A^\top)$.

**Proof.** Items 1 and 2 are easy to verify.

For item 3: ($\Rightarrow$) If $Pv = 0$, then $v = A^\top (AA^\top)^{-1} Av$. Letting $\lambda = (AA^\top)^{-1} Av$ yields $v = A^\top \lambda$.

($\Leftarrow$) Suppose $v = A^\top \lambda$, then

$$Pv = (I - A^\top (AA^\top)^{-1} A) A^\top \lambda = A^\top \lambda - A^\top \lambda = 0.$$

Similar to the derivation of $P$, we can obtain the projection onto $\Omega$:

$$\text{minimize} \quad \frac{1}{2}\|x - v\|^2$$
$$\text{subject to} \quad Ax = b$$

(Write down the Lagrange function and KKT condition, and solve for $(x, \lambda)$.)

The projection $\Pi$ of $v$ onto $\Omega$ is

$$\Pi(v) = Pv - A^\top (AA^\top)^{-1}b$$

**Proposition.** Let $\boldsymbol{x}^* \in \mathbb{R}^n$ be feasible (i.e., $\boldsymbol{A}\boldsymbol{x}^* = \boldsymbol{b}$), then $\boldsymbol{P}\nabla f(\boldsymbol{x}^*) = \boldsymbol{0}$ iff $\boldsymbol{x}^*$ satisfies the Lagrange condition.

**Proof.** We have

$$
\begin{aligned}
\boldsymbol{P}\nabla f(\boldsymbol{x}^*) = \boldsymbol{0} \quad &\Longleftrightarrow \quad \nabla f(\boldsymbol{x}^*) \in \mathcal{N}(\boldsymbol{P}) \\
&\Longleftrightarrow \quad \nabla f(\boldsymbol{x}^*) \in \mathcal{R}(\boldsymbol{A}^\top) \\
&\Longleftrightarrow \quad \nabla f(\boldsymbol{x}^*) = -\boldsymbol{A}^\top \boldsymbol{\lambda}^* \text{ for some } \boldsymbol{\lambda}^* \in \mathbb{R}^m
\end{aligned}
$$

Now we are ready to write down explicitly the PG:

$$
\begin{aligned}
\boldsymbol{x}^{(k+1)} &= \Pi(\boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)})) && (\because \text{PG definition}) \\
&= \boldsymbol{P}(\boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)})) - \boldsymbol{A}^\top (\boldsymbol{A}\boldsymbol{A}^\top)^{-1}\boldsymbol{b} && (\because \text{Relation of } \Pi \text{ and } \boldsymbol{P}) \\
&= \boldsymbol{P}\boldsymbol{x}^{(k)} - \boldsymbol{A}^\top (\boldsymbol{A}\boldsymbol{A}^\top)^{-1}\boldsymbol{b} - \boldsymbol{P}\alpha_k \nabla f(\boldsymbol{x}^{(k)}) \\
&= \Pi(\boldsymbol{x}^{(k)}) - \alpha_k \boldsymbol{P}\nabla f(\boldsymbol{x}^{(k)}) && (\because \text{Relation of } \Pi \text{ and } \boldsymbol{P}) \\
&= \boldsymbol{x}^{(k)} - \alpha_k \boldsymbol{P}\nabla f(\boldsymbol{x}^{(k)}) && (\because \boldsymbol{x}^{(k)} \in \Omega)
\end{aligned}
$$

The only difference from standard gradient method is the additional $\boldsymbol{P}$.

Note that if $\boldsymbol{x}^{(0)} \in \Omega$, then $\boldsymbol{x}^{(k)} \in \Omega$ for all $k$.

Now we can consider the choice of $\alpha_k$. For example, we can use the projected steepest descent (PSD) method:

$$\alpha_k = \operatorname*{arg\,min}_{\alpha > 0} f(\boldsymbol{x}^{(k)} - \alpha \boldsymbol{P} \nabla f(\boldsymbol{x}^{(k)}))$$

**Theorem.** Let $x^{(k)}$ be generated by PSD. If $P\nabla f(x^{(k)}) \neq 0$, then $f(x^{(k+1)}) < f(x^{(k)})$.

**Proof.** For such $x^{(k)}$, consider the line search function

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

Then we have

$$\phi'(\alpha) = -\nabla f(x^{(k)} - \alpha P\nabla f(x^{(k)}))^\top P\nabla f(x^{(k)}).$$

Hence

$$\phi'(0) = -\nabla f(x^{(k)})^\top P\nabla f(x^{(k)})$$
$$= -\nabla f(x^{(k)})^\top P^2 \nabla f(x^{(k)})$$
$$= -\|P\nabla f(x^{(k)})\|^2 < 0,$$

and therefore $\phi(\alpha_k) < \phi(0)$, i.e., $f(x^{(k+1)}) < f(x^{(k)})$.

$P\nabla f(x^*) = 0$ is sufficient for global optimality if $f$ is convex:

**Theorem.** Let $f$ be convex and $x^*$ be feasible. Then $P\nabla f(x^*) = 0$ iff $x^*$ is a global minimizer.

**Proof.** From the previous proposition and convexity of $f$, we know

$$P\nabla f(x^*) = 0 \quad \Longleftrightarrow \quad x^* \text{ satisfies the Lagrange condition}$$
$$\Longleftrightarrow \quad x^* \text{ is a global minimizer}$$

**Lagrange algorithm**

We first consider the Lagrange algorithm for equality-constrained optimization:

$$\text{minimize} \quad f(x)$$
$$\text{subject to} \quad h(x) = 0$$

where $f, h \in \mathcal{C}^2$.

Recall the Lagrange function $l : \mathbb{R}^{n+m} \to \mathbb{R}$ is

$$l(x, \lambda) = f(x) + h(x)^\top \lambda.$$

We denote its Hessian with respect to $x$ by

$$\nabla_x^2 l(x, \lambda) = \nabla_x^2 f(x) + D_x^2 h(x)^\top \lambda \in \mathbb{R}^{n \times n}$$

Recall the Lagrange condition is

$$\nabla f(\boldsymbol{x}) + D\boldsymbol{h}(\boldsymbol{x})^\top \boldsymbol{\lambda} = \boldsymbol{0} \in \mathbb{R}^n$$
$$\boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{0} \in \mathbb{R}^m$$

The **Lagrange algorithm** is given by

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k(\nabla f(\boldsymbol{x}^{(k)}) + D\boldsymbol{h}(\boldsymbol{x}^{(k)})^\top \boldsymbol{\lambda}^{(k)})$$
$$\boldsymbol{\lambda}^{(k+1)} = \boldsymbol{\lambda}^{(k)} + \beta_k \boldsymbol{h}(\boldsymbol{x}^{(k)})$$

which is like "gradient descent for $\boldsymbol{x}$" and "gradient ascent for $\boldsymbol{\lambda}$" of $l$.

Here $\alpha_k, \beta_k \geq 0$ are step sizes. WLOG, we can assume $\alpha_k = \beta_k$ for all $k$ by scaling $\boldsymbol{\lambda}^{(k)}$ properly.

It is easy to verify that, if $(\boldsymbol{x}^{(k)}, \boldsymbol{\lambda}^{(k)}) \to (\boldsymbol{x}^*, \boldsymbol{\lambda}^*)$, then $(\boldsymbol{x}^*, \boldsymbol{\lambda}^*)$ satisfies the Lagrange condition.

We denote $w = [x; \lambda] \in \mathbb{R}^{n+m}$ and

$$u(w) = \begin{bmatrix} x - \alpha(\nabla f(x) + Dh(x)^\top \lambda) \\ \lambda + \alpha h(x) \end{bmatrix} \in \mathbb{R}^{n+m}$$

Hence the Jacobian of $u$ is

$$\nabla u(w) = I + \alpha \begin{bmatrix} -\nabla_x^2 l(x, \lambda) & -Dh(x)^\top \\ Dh(x) & 0 \end{bmatrix} \in \mathbb{R}^{(n+m) \times (n+m)}$$

Note that

$$w^* = [x^*; \lambda^*] \text{ is a KKT point} \iff w^* = u(w^*)$$

We denote

$$M := \begin{bmatrix} -\nabla_x^2 l(x^*, \lambda^*) & -Dh(x^*)^\top \\ Dh(x^*) & 0 \end{bmatrix}$$

and hence $\nabla u(w^*) = I + \alpha M$.

Now we study the (local) convergence of the Lagrange algorithm when $x^*$ is a regular point and $\nabla_x^2 l(x^*, \lambda^*) \succ 0$. For simplicity, we assume $\alpha_k = \alpha$ (constant step size).

**Claim 1.** $\|\nabla u(w^*)\| < 1$ if $\alpha > 0$ is sufficiently small.

**Proof (Claim 1).** It suffices to show real part of any eigenvalue of $M$ is $< 0$.

Let $\lambda$ be an eigenvalue of $M$ and $w = [x; \lambda] \in \mathbb{C}^{n+m}$ be a corresponding eigenvector, i.e., $Mw = \lambda w$. (Note $w \neq 0$.)

If $x = 0$, then

$$Mw = \begin{bmatrix} -\nabla_x^2 l(x^*, \lambda^*) & -Dh(x^*)^\top \\ Dh(x^*) & 0 \end{bmatrix} \begin{bmatrix} 0 \\ \lambda \end{bmatrix} = \begin{bmatrix} -Dh(x^*)^\top \lambda \\ 0 \end{bmatrix} = \lambda \begin{bmatrix} 0 \\ \lambda \end{bmatrix}$$

But $Dh(x^*)$ has full row rank, so $\lambda = 0$, and hence $w = 0$, contradiction.

**Proof (Claim 1) cont.** Therefore $x \neq 0$. We know [*]

$$\Re(w^H M w) = \Re(w^H \lambda w) = \Re(\lambda)\|w\|^2$$

On the other hand [†]

$$\Re(w^H M w) = -\Re(x^H \nabla_x^2 l(x^*, \lambda^*)x) < 0$$

Equating the two yields $\Re(\lambda) < 0$.

As all eigenvalues of $M$ have negative real part, we know $\|I + \alpha M\| < 1$ for sufficiently small $\alpha > 0$.

This completes the proof of Claim 1.

[*] $w^H$ is the complex conjugate of $w$.
[†] Recall that if $Q \succ 0$, then $x^H Q x = \|\Re(x)\|_Q^2 + \|\Im(x)\|_Q^2$.

**Claim 2.** There exist $\eta > 0$ and $\kappa \in (0, 1)$ such that

$$\|\nabla u(w)\| \leq \kappa < 1, \quad \forall\, w \in B(w^*, \eta)$$

where $B(w^*, \eta) = \{w : \|w - w^*\| \leq \eta\}$.

**Proof (Claim 2).** The claim follows $\|\nabla u(w^*)\| < 1$ in Claim 1 and the continuity of $\nabla u$.

**Claim 3.** If $w^{(0)} \in B(w^*, \eta)$, then for all $k$ there is

$$\|w^{(k+1)} - w^*\| \leq \kappa \|w^{(k)} - w^*\|$$

**Proof (Claim 3).** Let $G : \mathbb{R}^{n+m} \to \mathbb{R}^{(n+m)\times(n+m)}$ be the function s.t.

$$u(w^{(k)}) - u(w^*) = G(w^{(k)})(w^{(k)} - w^*)$$

from the Mean Value Theorem. Hence

$$
\begin{aligned}
\|w^{(k+1)} - w^*\| &= \|u(w^{(k)}) - u(w^*)\| \\
&= \|G(w^{(k)})(w^{(k)} - w^*)\| \\
&\leq \|G(w^{(k)})\| \cdot \|w^{(k)} - w^*\| \\
&\leq \kappa \|w^{(k)} - w^*\|
\end{aligned}
$$

Claim 3 implies that locally $w^{(k)} \to w^*$ at a linear rate.

Now consider **Lagrange algorithm** for inequality-constrained optimization:

$$\text{minimize} \quad f(x)$$
$$\text{subject to} \quad g(x) \leq 0$$

The Lagrange function is

$$l(x, \mu) = f(x) + g(x)^\top \mu$$

The Lagrange condition is

$$\nabla f(x) + Dg(x)^\top \mu = 0$$
$$g(x) \leq 0$$
$$\mu \geq 0$$
$$g(x)^\top \mu = 0$$

The Lagrange algorithm is given by

$$x^{(k+1)} = x^{(k)} - \alpha_k(\nabla f(x^{(k)}) + Dg(x^{(k)})^\top \mu^{(k)})$$
$$\mu^{(k+1)} = [\mu^{(k)} + \beta_k g(x^{(k)})]_+$$

where $[\cdot]_+$ means $\max(\cdot, 0)$ componentwisely.

We denote $w = [x; \mu] \in \mathbb{R}^{n+p}$ and

$$\Pi(w) = \begin{bmatrix} x \\ [\mu]_+ \end{bmatrix}, \qquad u(w) = \begin{bmatrix} x - \alpha(\nabla f(x) + Dg(x)^\top \mu) \\ \mu + \alpha g(x) \end{bmatrix}$$

It is easy to verify that

$$w^* = [x^*; \mu^*] \text{ is a KKT point} \quad \Longleftrightarrow \quad w^* = \Pi(u(w^*))$$

Let $w^*$ be a KKT point, and

$$g(w^*) = \begin{bmatrix} g_A(w^*) \\ g_I(w^*) \end{bmatrix} \in \mathbb{R}^p = \mathbb{R}^{p_1+p_2}, \quad \text{where} \quad \begin{matrix} 0 = g_A(w) & \in \mathbb{R}^{p_1} \\ 0 < g_I(w) & \in \mathbb{R}^{p_2} \end{matrix}$$

"$A$" and "$I$" stand for "active" and "inactive".

Similarly, denote

$$\mu = \begin{bmatrix} \mu_A \\ \mu_I \end{bmatrix}, \quad w_A = \begin{bmatrix} x \\ \mu_A \end{bmatrix}, \quad u_A(w_A) = \begin{bmatrix} x - \alpha(\nabla f(x) + Dg_A(x)^\top \mu_A) \\ \mu_A + \alpha g_A(x) \end{bmatrix}$$

and hence

$$\nabla u_A(w_A) = I + \alpha \begin{bmatrix} -\nabla_x^2 l(x, \mu_A) & -Dg_A(x)^\top \\ Dg_A(x) & 0 \end{bmatrix} \in \mathbb{R}^{(n+p_1) \times (n+p_1)}$$

Now we study the (local) convergence of the Lagrange algorithm when $x^*$ is a regular point and $\nabla_x^2 l(x^*, \lambda^*) \succ 0$. For simplicity, we assume $\alpha_k = \alpha$ (constant step size).

We again define $G$ such that

$$u(w^{(k)}) - u(w^*) = G(w^{(k)})(w^{(k)} - w^*)$$

using Mean Value Theorem. Let

$$M = \begin{bmatrix} -\nabla_x^2 l(x^*, \mu_A^*) & -Dg_A(x^*)^\top \\ Dg_A(x^*) & 0 \end{bmatrix} \in \mathbb{R}^{(n+p_1) \times (n+p_1)}$$

Similar as before, we can show all eigenvalues of $M$ have negative real part, and hence $\|I + \alpha M\| < 1$ for $\alpha$ sufficiently small.

Also note that $\mu_I^* = 0$ as it corresponds to the inactive constraints.

**Claim 1.** There exist $\eta > 0$ and $\kappa_A \in (0, 1)$, such that

$$\|\nabla \boldsymbol{u}_A(\boldsymbol{w}_A)\| \leq \kappa_A$$
$$\boldsymbol{g}_I(\boldsymbol{x}) \leq -\delta e$$

for all $\boldsymbol{w} \in B(\boldsymbol{w}^*, \eta)$.

**Proof.** Note that $\boldsymbol{g}_I(\boldsymbol{w}^*) < \boldsymbol{0}$. Others are similar as before.

Now we set the following values:

- Let $\kappa = \max\{1, \|\boldsymbol{G}(\boldsymbol{w})\| : \boldsymbol{w} \in B(\boldsymbol{w}^*, \eta)\} \geq 1$

- Let $\varepsilon > 0$ be small enough such that $\varepsilon \kappa^{\varepsilon/(\alpha\delta)} \leq \eta$.

- Let $k_0 = \lceil \varepsilon/(\alpha\delta) \rceil$.

- Let $\boldsymbol{w}^{(0)} \in B(\boldsymbol{w}^*, \varepsilon)$.

**Claim 2.** For any $k \leq k_0$, there is $\|w^{(k)} - w^*\| \leq \varepsilon \kappa^k$.

**Proof (Claim 2).** We use induction.

First, there is $\|w^{(0)} - w^*\| \leq \varepsilon = \varepsilon \kappa^0$.

Assume the claim holds for $k$, then

$$
\begin{aligned}
\|w^{(k+1)} - w^*\| &\leq \|G(w^{(k)})\| \cdot \|w^{(k)} - w^*\| \\
&\leq \kappa \cdot \|w^{(k)} - w^*\| \\
&\leq \kappa \cdot (\varepsilon \kappa^k) \\
&= \varepsilon \kappa^{k+1}
\end{aligned}
$$

which completes the proof of the claim.

From Claim 2, we know $\|w^{(k)} - w^*\| \leq \eta$ for $k = 0, \ldots, k_0$.

**Claim 3.** There is $\mu_I^{(0)} \geq \cdots \geq \mu_I^{(k_0)} = 0$.

**Proof (Claim 3).** We know $g_I(x^{(k)}) \leq -\delta e$ for $k = 0, \ldots, k_0$. Also

$$\mu_I^{(k+1)} = [\mu_I^{(k)} + \alpha g_I(x^{(k)})]_+ \leq [\mu_I^{(k)} - \alpha \delta e]_+ \leq \mu_I^{(k)}$$

which implies that $\mu_I^{(k)}$ is non-increasing.

Suppose $\mu_i^{(k_0)} > 0$ for some $i \in I$ (index set of inactive constraints), then

$$0 < \mu_i^{(k_0)} = \mu_i^{(k_0-1)} + \alpha g_i(x^{(k_0-1)}) = \cdots$$

$$= \mu_i^{(0)} + \alpha \sum_{k=0}^{k_0-1} g_i(x^{(k)}) \leq \mu_i^{(0)} - \alpha \delta k_0 \leq \varepsilon - \alpha \delta k_0$$

since $\mu_i^{(0)} \leq \|w^{(0)} - w^*\| \leq \varepsilon$. But this contradicts to $k_0 = \lceil \frac{\epsilon}{\alpha \delta} \rceil \geq \frac{\epsilon}{\alpha \delta}$.

Therefore, within $k_0$ iterations, $\mu_I^{(k)} = 0$.

**Claim 4.** For any $k \geq k_0$, there are

$$\mu_I^{(k)} = 0$$
$$\|w^{(k)} - w^*\| \leq \eta$$
$$\|w_A^{(k+1)} - w_A^*\| \leq \kappa_A \|w_A^{(k)} - w_A^*\|$$

**Proof (Claim 4).** The first two hold for $k = k_0$ (by Claims 3 & 2 resp.), and

$$
\begin{aligned}
\|w_A^{(k_0+1)} - w_A^*\| &= \|\,\Pi(u_A(w_A^{(k_0)})) - \Pi(u_A(w_A^*))\| \\
&\leq \|u_A(w_A^{(k_0)}) - u_A(w_A^*)\| \\
&\leq \|G_A(w_A^{(k)})\| \cdot \|w_A^{(k_0)} - w_A^*\| \\
&\leq \kappa_A \cdot \|w_A^{(k_0)} - w_A^*\|
\end{aligned}
$$

**Proof (Claim 4) cont.**

Assume the results hold for $k \geq k_0$, then from $g_I(w^{(k)}) \leq -\delta e$, we have

$$\mu_I^{(k+1)} = [\mu_I^{(k)} + \alpha g_I(x^{(k)})]_+ \leq [0 - \alpha\delta e]_+ = 0$$

Note that this implies $\|w_A^{(k+1)} - w_A^*\| = \|w^{(k+1)} - w^*\|$ for all $k \geq k_0$.

Moreover, we have $w_A^{(k+2)} = \Pi(u_A(w_A^{(k+1)}))$ and

$$\|w_A^{(k+2)} - w_A^*\| \leq \kappa_A \cdot \|w_A^{(k+1)} - w_A^*\| \leq \eta$$

which completes the proof.

**Remark.** Claim 4 implies that locally $w^{(k)} \to w^*$ at a linear rate: if $w^{(0)}$ is sufficiently close to $w^*$, then $w^{(k)} \to w^*$ linearly, provided that $x^*$ is a regular KKT point and $\nabla_x^2 l(x^*, \lambda^*) \succ 0$.

## Penalty method

We consider constrained optimization

$$\text{minimize} \quad f(x)$$
$$\text{subject to} \quad x \in \Omega$$

Note that such problem conceptually include optimization problems with equality and inequality constraints. For example, $\Omega = \{x \in \mathbb{R}^n : g(x) \le 0\}$.

Instead of the constrained problem, we consider to impose penalty if $x \in \Omega$ is violated:

$$\text{minimize} \ f(x) + \gamma P(x)$$

where $P : \mathbb{R}^n \to \mathbb{R}_+$ is the penalty function, and $\gamma > 0$ is the penalty (weight) parameter.

**Definition.** The function $P : \mathbb{R}^n \to \mathbb{R}_+$ is called a **penalty function** if

1. $P$ is continuous.

2. $P(x) \geq 0$ for all $x$.

3. $P(x) = 0$ iff $x \in \Omega$.

**Example.** Let $\Omega = \{x \in \mathbb{R}^n : g(x) \leq 0 \in \mathbb{R}^p\}$, then we can choose

$$P(x) = \sum_{i=1}^{p} [g_i(x)]_+$$

$$P(x) = \sum_{i=1}^{p} ([g_i(x)]_+)^2$$

and so on.

**Example.** Let $g(x) = [g_1(x); g_2(x)]$ where $g_1(x) = x - 2$ and $g_2(x) = -(x+1)^3$. Consider the constraint set

$$\Omega = \{x \in \mathbb{R} : g_1(x) \leq 0, \ g_2(x) \leq 0\}$$

Then we have

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

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

We can set

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

**Example.** Consider the problem below with $Q \succ 0$:

$$\text{minimize} \quad x^\top Q x$$
$$\text{subject to} \quad \|x\|^2 = 1$$

We can set the penalty function $P(x) = (\|x\|^2 - 1)^2$ (which is differentiable), and consider

$$\text{minimize} \quad x^\top Q x + \gamma(\|x\|^2 - 1)^2$$

For any fixed $\gamma > 0$, the FONC of its solution $x_\gamma$ is

$$2Q x_\gamma + 4\gamma(\|x_\gamma\|^2 - 1)x_\gamma = 0$$

which yields

$$Q x_\gamma = 2\gamma(1 - \|x_\gamma\|^2)x_\gamma = \lambda_\gamma x_\gamma$$

where $\lambda_\gamma := 2\gamma(1 - \|x_\gamma\|^2)$ is a scalar. This means $\lambda_\gamma \in (0, \lambda_{\max}]$ is an eigenvalue of $Q$, and $x_\gamma$ is a corresponding eigenvector. Note that

$$0 < 1 - \|x_\gamma\|^2 \leq \frac{\lambda_{\max}}{2\gamma} = \mathcal{O}\left(\frac{1}{\gamma}\right).$$

We have converted constrained problem into unconstrained ones. Now define

$$q(\gamma_k, \boldsymbol{x}) = f(\boldsymbol{x}) + \gamma_k P(\boldsymbol{x})$$
$$\boldsymbol{x}^{(k)} = \arg\min_{\boldsymbol{x} \in \mathbb{R}^n} q(\gamma_k, \boldsymbol{x})$$

for every $k \in \mathbb{N}$.

The idea is to let $\gamma_k$ increase (hence greater penalty) and apply an unconstrained optimization method to solve for $\boldsymbol{x}^{(k)}$ for each $k$.

Then we hope that an accumulation point[‡] of $\{\boldsymbol{x}^{(k)}\}$ is a KKT point $\boldsymbol{x}^*$.

[‡]$\boldsymbol{x}^*$ is called an *accumulation point* (also called *limit point*) of $\{\boldsymbol{x}^{(k)}\}$ if there exists a subsequence of $x^{(k)}$ that converges to $\boldsymbol{x}^*$.

Now let $\gamma_k > 0$ be increasing, we have a series of claims.

**Claim 1.** $q(\gamma_k, \boldsymbol{x}^{(k)}) \leq q(\gamma_{k+1}, \boldsymbol{x}^{(k+1)})$.

**Proof (Claim 1).** Since $\boldsymbol{x}^{(k)}$ is optimal to $q(\gamma_k, \boldsymbol{x})$, we know

$$q(\gamma_k, \boldsymbol{x}^{(k)}) \leq q(\gamma_k, \boldsymbol{x}^{(k+1)})$$

Furthermore, since $\gamma_k < \gamma_{k+1}$, we know

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

Combining the two verifies the claim.

**Claim 2.** $P(\boldsymbol{x}^{(k+1)}) \leq P(\boldsymbol{x}^{(k)})$.

**Proof (Claim 2).** By the optimality of $\boldsymbol{x}^{(k)}$ and $\boldsymbol{x}^{(k+1)}$ for their own problems, we know

$$q(\gamma_k, \boldsymbol{x}^{(k)}) \leq q(\gamma_k, \boldsymbol{x}^{(k+1)})$$
$$q(\gamma_{k+1}, \boldsymbol{x}^{(k+1)}) \leq q(\gamma_{k+1}, \boldsymbol{x}^{(k)})$$

which are

$$f(\boldsymbol{x}^{(k)}) + \gamma_k P(\boldsymbol{x}^{(k)}) \leq f(\boldsymbol{x}^{(k+1)}) + \gamma_k P(\boldsymbol{x}^{(k+1)})$$
$$f(\boldsymbol{x}^{(k+1)}) + \gamma_{k+1} P(\boldsymbol{x}^{(k+1)}) \leq f(\boldsymbol{x}^{(k)}) + \gamma_{k+1} P(\boldsymbol{x}^{(k)})$$

Adding the two above yields

$$(\gamma_{k+1} - \gamma_k) P(\boldsymbol{x}^{(k+1)}) \leq (\gamma_{k+1} - \gamma_k) P(\boldsymbol{x}^{(k)})$$

Recalling $\gamma_{k+1} - \gamma_k > 0$ completes the proof.

**Claim 3.** $f(\boldsymbol{x}^{(k+1)}) \geq f(\boldsymbol{x}^{(k)})$.

**Proof (Claim 3).** Since $q(\gamma_k, \boldsymbol{x}^{(k)}) \leq q(\gamma_k, \boldsymbol{x}^{(k+1)})$, we know

$$f(\boldsymbol{x}^{(k)}) + \gamma_k P(\boldsymbol{x}^{(k)}) \leq f(\boldsymbol{x}^{(k+1)}) + \gamma_k P(\boldsymbol{x}^{(k+1)})$$

From Claim 2, we know $P(\boldsymbol{x}^{(k+1)}) \leq P(\boldsymbol{x}^{(k)})$, hence

$$f(\boldsymbol{x}^{(k+1)}) \geq f(\boldsymbol{x}^{(k)}) + \gamma_k(P(\boldsymbol{x}^{(k)}) - P(\boldsymbol{x}^{(k+1)})) \geq f(\boldsymbol{x}^{(k)})$$

**Claim 4.** $f(\boldsymbol{x}^*) \geq q(\gamma_k, \boldsymbol{x}^{(k)}) \geq f(\boldsymbol{x}^{(k)})$.

**Proof (Claim 4).** We know $P(\boldsymbol{x}^*) = 0$, and hence

$$
\begin{aligned}
f(\boldsymbol{x}^*) &= q(\gamma_k, \boldsymbol{x}^*) \\
&\geq q(\gamma_k, \boldsymbol{x}^{(k)}) \\
&= f(\boldsymbol{x}^{(k)}) + \gamma_k P(\boldsymbol{x}^{(k)}) \\
&\geq f(\boldsymbol{x}^{(k)})
\end{aligned}
$$

**Theorem.** Suppose $f$ is continuous and $\gamma_k \uparrow \infty$. Then any accumulation point of $\{x^{(k)}\}$ is a solution to the constrained problem.

**Proof.** For simplicity, let $x^{(k)}$ denote the subsequence which converges to $\widehat{x}$.

Since $f(x^{(k)}) \leq f(x^*)$ for all $k$ (by Claim 4), we know

$$f(x^*) \geq \lim_{k \to \infty} f(x^{(k)}) = f(\widehat{x})$$

Note that $q(\gamma_k, x^{(k)})$ is nondecreasing in $k$ (by Claim 1) and bounded above by $f(x^*)$ (by Claim 4), we know $q(\gamma_k, x^{(k)}) \uparrow q^*$ for some $q^* \in \mathbb{R}$. Hence,

$$\gamma_k P(x^{(k)}) = q(\gamma_k, x^{(k)}) - f(x^{(k)}) \to q^* - f(\widehat{x})$$

Since $\gamma_k \to \infty$, we know $P(x^{(k)}) \to 0$. Since $P$ is continuous, we know $P(\widehat{x}) = 0$, i.e., $\widehat{x}$ is feasible. Therefore $\widehat{x}$ is optimal since $f(\widehat{x}) \leq f(x^*)$.

Penalty method requires solving one instance of

$$\text{minimize } f(x) + \gamma P(x)$$

with $\gamma = \gamma_k$ for every $k$.

Is it possible to obtain the solution with a single $\gamma$?

**Definition.** We call $P$ an **exact penalty** if there exists $\gamma > 0$ such that the solution $x^*$ of the unconstrained problem

$$\text{minimize } f(x) + \gamma P(x)$$

is also a solution of the constrained problem

$$\text{minimize} \quad f(x)$$
$$\text{subject to} \quad x \in \Omega$$

However it turns out that it may be necessary for an exact penalty $P$ to be non-differentiable.

**Proposition.** Let $\Omega$ be convex, $x^*$ is on the boundary of $\Omega$. If there exists a feasible direction $d$ at $x^*$ such that $d^\top \nabla f(x^*) > 0$, then an exact penalty $P$ must be non-differentiable.

**Proof.** Suppose not, then $\nabla P(x^*) = 0$ since $P(x) = 0$ for all $x \in \Omega$. Let $g(x) = f(x) + \gamma P(x)$, then

$$\nabla g(x^*) = \nabla f(x^*) + \gamma \nabla P(x^*) = \nabla f(x^*)$$

and hence $d^\top g(x^*) = d^\top \nabla f(x^*) > 0$, which means $x^*$ is not a local minimizer of $g$, contradiction.

**Example.** Consider the problem

$$\text{minimize} \quad 5 - 3x$$
$$\text{subject to} \quad x \in [0, 1]$$

We can see $x^* = 1$ which is on the boundary, and $f'(x^*) = -3$ aligns with the feasible direction $d = -1$ at $x^*$.

If we use a differentiable penalty function $P$, then $P'(x^*) = 0$. Let

$$g(x) = f(x) + \gamma P(x),$$

then $g'(x^*) = f'(x^*) + \gamma P'(x^*) = -3 \neq 0$, which means $P$ cannot be an exact penalty function.

**Remark.** However, if $\boldsymbol{d}^\top \nabla f(\boldsymbol{x}^*) \leq 0$ for any feasible direction $\boldsymbol{d}$ at $\boldsymbol{x}$, we may still be able to find a differentiable exact penalty function $P$.