

Lecture 9: Penalty Methods and Multiplier Methods

- Quadratic Penalty Methods
- Introduction to Multiplier Methods

Two Convergence Mechanisms

Consider the equality constrained problem

$$\begin{array}{ll}\text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in X, \mathbf{h}(\mathbf{x}) = \mathbf{0},\end{array}$$

where $f : \mathbb{R}^n \mapsto \mathbb{R}$ and $\mathbf{h} : \mathbb{R}^n \mapsto \mathbb{R}^m$ are continuous, and X is closed.

- The quadratic penalty method:

$$\mathbf{x}^r = \arg \min_{\mathbf{x} \in X} L_{c^r}(\mathbf{x}, \boldsymbol{\lambda}^r) \equiv f(\mathbf{x}) + (\boldsymbol{\lambda}^r)' \mathbf{h}(\mathbf{x}) + \frac{c^r}{2} \|\mathbf{h}(\mathbf{x})\|^2$$

where the $\boldsymbol{\lambda}^r$ is a bounded sequence and c^r satisfies $0 < c^r < c^{r+1}$ for all r and $c^r \rightarrow \infty$.

- Mechanism 1 for convergence: taking $\boldsymbol{\lambda}^r$ close to a Lagrange multiplier vector
 - ★ Assume $X = \mathbb{R}^n$ and $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ is a local min-Lagrange multiplier pair satisfying the 2nd order sufficiency conditions

- ★ For c sufficiently large, \mathbf{x}^* is a strict local min of $L_c(\cdot, \boldsymbol{\lambda}^*)$
- Mechanism 2 for convergence: Taking c^r very large
 - ★ For large c and any $\boldsymbol{\lambda}$, we have

$$L_c(\cdot, \boldsymbol{\lambda}) \approx \begin{cases} f(\mathbf{x}) & \text{if } \mathbf{x} \in X \text{ and } h(\mathbf{x}) = \mathbf{0} \\ \infty & \text{otherwise} \end{cases}$$

- Example:

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) = \frac{1}{2}(x_1^2 + x_2^2) \\ \text{subject to} & x_1 = 1 \end{array}$$

We have $\mathbf{x}^* = (1, 0)$, $\lambda^* = -1$ and

$$\begin{aligned} L_c(\mathbf{x}, \lambda) &= \frac{1}{2}(x_1^2 + x_2^2) + \lambda(x_1 - 1) + \frac{c}{2}(x_1 - 1)^2 \\ x_1(\lambda, c) &= \frac{c - \lambda}{c + 1}, \quad x_2(\lambda, c) = 0 \end{aligned}$$

Global Convergence

- Suppose $c^r \rightarrow \infty$. Then every limit point of $\{\mathbf{x}^r\}$ is a global min.
- **Proof:** The optimal value of the problem is

$$f^* = \inf_{\mathbf{h}(\mathbf{x})=\mathbf{0}, \mathbf{x} \in X} L_{c^r}(\mathbf{x}, \boldsymbol{\lambda}^r).$$

We have $L_{c^r}(\mathbf{x}^r, \boldsymbol{\lambda}^r) \leq L_{c^r}(\mathbf{x}, \boldsymbol{\lambda}^r)$, $\forall \mathbf{x} \in X$ so taking the inf of the RHS over $\mathbf{x} \in X$, $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ yields

$$L_{c^r}(\mathbf{x}^r, \boldsymbol{\lambda}^r) = f(\mathbf{x}^r) + (\boldsymbol{\lambda}^r)' \mathbf{h}(\mathbf{x}^r) + \frac{c^r}{2} \|\mathbf{h}(\mathbf{x}^r)\|^2 \leq f^*$$

Let $(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}})$ be a limit point of $\{\mathbf{x}^r, \boldsymbol{\lambda}^r\}$. Without loss of generality, assume that $\{\mathbf{x}^r, \boldsymbol{\lambda}^r\} \rightarrow (\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}})$. Taking the limsup above

$$f(\bar{\mathbf{x}}) + \bar{\boldsymbol{\lambda}}' \mathbf{h}(\bar{\mathbf{x}}) + \limsup_{r \rightarrow \infty} \frac{c^r}{2} \|\mathbf{h}(\mathbf{x}^r)\|^2 \leq f^*$$

By $\|\mathbf{h}(\mathbf{x}^r)\|^2 \geq 0$ and $\{c^r\} \rightarrow \infty$, we have $\mathbf{h}(\mathbf{x}^r) \rightarrow \mathbf{0}$ and $\mathbf{h}(\bar{\mathbf{x}}) = \mathbf{0}$. Hence, $\bar{\mathbf{x}}$ is feasible, and since the above inequality implies $f(\bar{\mathbf{x}}) \leq f^*$, so $\bar{\mathbf{x}}$ is optimal.

Lagrange Multiplier Estimates

- Assume that $X = \mathbb{R}^n$, and f and h are continuously differentiable. Let $\{\lambda^r\}$ be bounded, and $\{c^r\} \rightarrow \infty$. Assume x^r satisfies $\nabla_x L_{c^r}(x^r, \lambda^r) = 0$ for all r , and that $x^r \rightarrow x^*$, where x^* is such that $\text{rank}(\nabla h(x^*)) = m$. Then $h(x^*) = 0$ and $\tilde{\lambda}^r \rightarrow \lambda^*$, where

$$\tilde{\lambda}^r = \lambda^r + c^r h(x^r), \quad \nabla_x L(x^*, \lambda^*) = 0.$$

- Proof:** We have

$$0 = \nabla_x L_{c^r}(x^r, \lambda^r) = \nabla f(x^r) + \nabla h(x^r) (\lambda^r + c^r h(x^r)) = \nabla f(x^r) + \nabla h(x^r) \tilde{\lambda}^r.$$

- Multiply with

$$(\nabla h(x^r)' \nabla h(x^r))^{-1} \nabla h(x^r)'$$

and take lim to obtain $\tilde{\lambda}^r \rightarrow \lambda^*$ with

$$\lambda^* = -(\nabla h(x^*)' \nabla h(x^*))^{-1} \nabla h(x^*)' \nabla f(x^*).$$

We also have $\nabla_x L(x^*, \lambda^*) = 0$ and $h(x^*) = 0$ (since $\tilde{\lambda}^r$ converges).

Practical Behaviors

- Three possibilities:
 - ★ The method breaks down because an \mathbf{x}^r with $\nabla_{\mathbf{x}} L_{c^r}(\mathbf{x}^r, \boldsymbol{\lambda}^r) \approx \mathbf{0}$ cannot be found.
 - ★ A sequence $\{\mathbf{x}^r\}$ with $\nabla_{\mathbf{x}} L_{c^r}(\mathbf{x}^r, \boldsymbol{\lambda}^r) \approx \mathbf{0}$ is obtained, but it either has no limit points, or for each of its limit points \mathbf{x}^* the matrix $\nabla \mathbf{h}(\mathbf{x}^*)$ has rank $< m$.
 - ★ A sequence $\{\mathbf{x}^r\}$ with $\nabla_{\mathbf{x}} L_{c^r}(\mathbf{x}^r, \boldsymbol{\lambda}^r) \approx \mathbf{0}$ is found and it has a limit point \mathbf{x}^* such that $\nabla \mathbf{h}(\mathbf{x}^*)$ has rank m . Then, \mathbf{x}^* together with $\boldsymbol{\lambda}^*$ [the corresponding limit point of $\{\boldsymbol{\lambda}^r + c^r \mathbf{h}(\mathbf{x}^r)\}$] satisfies the first-order necessary conditions.
- Ill-conditioning: The condition number of the Hessian $\nabla_{\mathbf{x}\mathbf{x}}^2 L_{c^r}(\mathbf{x}^r, \boldsymbol{\lambda}^r)$ tends to increase with c^r .
- To overcome ill-conditioning:

- ★ Use Newton-like method (and double precision).
- ★ Use good starting points.
- ★ Increase $\{c^r\}$ at a moderate rate (if $\{c^r\}$ is increased at a fast rate, $\{x^r\}$ converges faster, but the likelihood of ill-conditioning is greater).

Inequality Constraints

- Convert them to equality constraints by using squared slack variables that are eliminated later.
- Convert inequality constraint $g_j(\mathbf{x}) \leq 0$ to equality constraint $g_j(\mathbf{x}) + z_j^2 = 0$.
- The penalty method solves problems of the form

$$\min_{\mathbf{x}, \mathbf{z}} \bar{L}_c(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = L_c(\mathbf{x}, \boldsymbol{\lambda}) + \sum_{j=1}^r \left[\mu_j (g_j(\mathbf{x}) + z_j^2) + \frac{c}{2} |g_j(\mathbf{x}) + z_j^2|^2 \right],$$

for various values of $\boldsymbol{\lambda}$, $\boldsymbol{\mu}$ and c .

- First minimize $\bar{L}_c(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ with respect to \mathbf{z} to compute $L_c(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ by

$$\min_{\mathbf{z}} \bar{L}_c(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = L_c(\mathbf{x}, \boldsymbol{\lambda}) + \sum_{j=1}^r \min_{z_j} \left[\mu_j (g_j(\mathbf{x}) + z_j^2) + \frac{c}{2} |g_j(\mathbf{x}) + z_j^2|^2 \right]$$

and then minimize $L_c(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ with respect to \mathbf{x} .

Multiplier Methods

- Recall that if $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ is a local min-Lagrange multiplier pair satisfying the 2nd order sufficiency conditions, then for c sufficiently large, \mathbf{x}^* is a strict local min of $L_c(\cdot, \boldsymbol{\lambda}^*)$.
- This suggests that for $\boldsymbol{\lambda}^r \approx \boldsymbol{\lambda}^*, \mathbf{x}^r \approx \mathbf{x}^*$.
- Hence it is a good idea to use $\boldsymbol{\lambda}^r \approx \boldsymbol{\lambda}^*$, such as

$$\boldsymbol{\lambda}^{r+1} = \tilde{\boldsymbol{\lambda}}^r = \boldsymbol{\lambda}^r + c^r \mathbf{h}(\mathbf{x}^r)$$

This is the (1st order) method of multipliers.

- Key advantages to be shown:
 - ★ Less ill-conditioning: It is not necessary that $c^r \rightarrow \infty$ (only that c^r exceeds some threshold).

★ Faster convergence when λ^r is updated than when λ^r is kept constant (whether $c^r \rightarrow \infty$ or not).

- Consider the equality constrained problem

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{h}(\mathbf{x}) = \mathbf{0}, \end{array}$$

where $f : \mathbb{R}^n \mapsto \mathbb{R}$ and $\mathbf{h} : \mathbb{R}^n \mapsto \mathbb{R}^m$ are continuously differentiable.

- The (1st order) multiplier method finds

$$\mathbf{x}^r = \arg \min_{\mathbf{x} \in \mathbb{R}^n} L_{c^r}(\mathbf{x}, \lambda^r) \equiv f(\mathbf{x}) + (\lambda^r)' \mathbf{h}(\mathbf{x}) + \frac{c^r}{2} \|\mathbf{h}(\mathbf{x})\|^2$$

and updates λ^r using

$$\lambda^{r+1} = \lambda^r + c^r \mathbf{h}(\mathbf{x}^r)$$

Convex Example

- Problem: $\min_{x_1=1} \frac{1}{2}(x_1^2 + x_2^2)$ with optimal solution $\mathbf{x}^* = (1, 0)$ and Lagrangian multiplier $\lambda^* = -1$.

- We have

$$\begin{aligned}\mathbf{x}^r &= \arg \min_{\mathbf{x} \in \mathbb{R}^n} L_{c^r}(\mathbf{x}, \lambda^r) = \left(\frac{c^r - \lambda^r}{c^r + 1}, 0 \right) \\ \lambda^{r+1} &= \lambda^r + c^r \left(\frac{c^r - \lambda^r}{c^r + 1} - 1 \right) \\ \lambda^{r+1} - \lambda^* &= \frac{\lambda^r - \lambda^*}{c^r + 1}\end{aligned}$$

- We see that:

- ★ $\lambda^r \rightarrow \lambda^* = -1$ and $\mathbf{x}^r \rightarrow \mathbf{x}^* = (1, 0)$ for every nondecreasing sequence $\{c^r\}$. It is NOT necessary to increase c^r to ∞ .
- ★ The convergence rate becomes faster as c^r becomes larger; in fact $\{|\lambda^r - \lambda^*|\}$ converges superlinearly if $c^r \rightarrow \infty$.

Nonconvex Example

- Problem: $\min_{x_1=1} \frac{1}{2}(-x_1^2 + x_2^2)$ with optimal solution $\mathbf{x}^* = (1, 0)$ and Lagrangian multiplier $\lambda^* = 1$.

- We have

$$\mathbf{x}^r = \arg \min_{\mathbf{x} \in \mathbb{R}^n} L_{c^r}(\mathbf{x}, \lambda^r) = \left(\frac{c^r - \lambda^r}{c^r - 1}, 0 \right)$$

provided $c^r > 1$ (otherwise the min does not exist)

$$\begin{aligned} \lambda^{r+1} &= \lambda^r + c^r \left(\frac{c^r - \lambda^r}{c^r - 1} - 1 \right) \\ \lambda^{r+1} - \lambda^* &= -\frac{\lambda^r - \lambda^*}{c^r - 1} \end{aligned}$$

- We see that:
 - ★ No need to increase c^r to ∞ for convergence; doing so results in faster convergence rate.
 - ★ To obtain convergence, c^r must eventually exceed the threshold 2.

Primal Functional

- Let $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ be a regular local min-Lagrangian pair satisfying the 2nd order sufficient conditions are satisfied.
- The primal functional

$$p(u) = \min_{\mathbf{h}(\mathbf{x})=u} f(\mathbf{x}),$$

defined for u in an open sphere centered at $u = 0$, and we have

$$p(0) = f(\mathbf{x}^*), \quad \nabla p(0) = -\boldsymbol{\lambda}^*.$$

- Two examples:

$$p(u) = \min_{x_1-1=u} \frac{1}{2}(x_1^2 + x_2^2) = \frac{1}{2}(1+u)^2, \quad p(0) = f(\mathbf{x}^*) = \frac{1}{2}, \quad p'(0) = 1 = -\lambda^*$$

and

$$p(u) = \min_{x_1-1=u} \frac{1}{2}(-x_1^2 + x_2^2) = -\frac{1}{2}(1+u)^2, \quad p'(0) = -1 = -\lambda^*$$

Augmented Lagrangian Minimization

- Break down the minimization of $L_c(\mathbf{x}, \boldsymbol{\lambda})$:

$$\begin{aligned}\min_{\mathbf{x}} L_c(\mathbf{x}, \boldsymbol{\lambda}) &= \min_{\mathbf{u}} \min_{\mathbf{h}(\mathbf{x})=\mathbf{u}} \left\{ f(\mathbf{x}) + \boldsymbol{\lambda}'\mathbf{h}(\mathbf{x}) + \frac{c}{2}\|\mathbf{h}(\mathbf{x})\|^2 \right\} \\ &= \min_{\mathbf{u}} \left\{ p(\mathbf{u}) + \boldsymbol{\lambda}'\mathbf{u} + \frac{c}{2}\|\mathbf{u}\|^2 \right\},\end{aligned}$$

where the minimization above is understood to be local in a neighborhood of $\mathbf{u} = \mathbf{0}$.

- Interpretation of this minimization: Penalized Primal Function $p(\mathbf{u}) + \frac{c}{2}\|\mathbf{u}\|^2$
- If c is sufficiently large, $p(\mathbf{u}) + \boldsymbol{\lambda}'\mathbf{u} + \frac{c}{2}\|\mathbf{u}\|^2$ is convex in a neighborhood of $\mathbf{0}$. Also, for $\boldsymbol{\lambda} \approx \boldsymbol{\lambda}^*$ and large c , the value $\min_{\mathbf{x}} L_c(\mathbf{x}, \boldsymbol{\lambda}) \approx p(\mathbf{0}) = f(\mathbf{x}^*)$.

Interpretation of This Method

- Geometric interpretation of the iteration

$$\lambda^{r+1} = \lambda^r + c^r h(x^r).$$

- If λ^r is sufficiently close to λ^* and/or c^r is sufficiently large, λ^{r+1} will be closer to λ^* than λ^r .
- c^r need not be increased to ∞ in order to obtain convergence; it is sufficient that c^r eventually exceeds some threshold level.
- If $p(u)$ is linear, convergence to λ^* will be achieved in one iteration.

Computational Aspects

- Key issue is how to select $\{c^r\}$.
- c^r should eventually become larger than the “threshold” of the given problem.
- c^0 should not be so large as to cause ill-conditioning at the 1st minimization.
- c^r should not be increased so fast that too much ill-conditioning is forced upon the unconstrained minimization too early.
- c^r should not be increased so slowly that the multiplier iteration has poor convergence rate.
- A good practical scheme is to choose a moderate value c^0 , and use $c^{r+1} = \beta c^r$, where $\beta > 1$ is a scalar (typically $\beta \in [5, 10]$ if a Newton like method is used).
- In practice the minimization of $L_{c^r}(\mathbf{x}, \boldsymbol{\lambda}^r)$ is typically inexact (usually exact asymptotically). In some variants of the method, only one Newton step per minimization is used (with safeguards).

Duality Framework

- Consider the problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) + \frac{c}{2} \|\mathbf{h}(\mathbf{x})\|^2 \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}^*\| < \epsilon, \quad \mathbf{h}(\mathbf{x}) = \mathbf{0}, \end{aligned}$$

where ϵ is small enough for a local analysis to hold based on the implicit function theorem, and c is large enough for the minimum to exist.

- Consider the dual function and its gradient

$$\begin{aligned} q_c(\boldsymbol{\lambda}) &= \min_{\|\mathbf{x} - \mathbf{x}^*\| < \epsilon} L_c(\mathbf{x}, \boldsymbol{\lambda}) = L_c(\mathbf{x}(\boldsymbol{\lambda}, c), \boldsymbol{\lambda}), \\ \nabla q_c(\boldsymbol{\lambda}) &= \nabla_{\boldsymbol{\lambda}} \mathbf{x}(\boldsymbol{\lambda}, c) \nabla_{\mathbf{x}} L_c(\mathbf{x}(\boldsymbol{\lambda}, c), \boldsymbol{\lambda}) + \mathbf{h}(\mathbf{x}(\boldsymbol{\lambda}, c)) = \mathbf{h}(\mathbf{x}(\boldsymbol{\lambda}, c)) \end{aligned}$$

We have $\nabla q_c(\boldsymbol{\lambda}^*) = \mathbf{h}(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 q_c(\boldsymbol{\lambda}^*) \succ 0$.

- The multiplier method is a steepest ascent iteration for maximizing q_{c^r}

$$\boldsymbol{\lambda}^{r+1} = \boldsymbol{\lambda}^r + c^r \nabla q_{c^r}(\boldsymbol{\lambda}^r).$$