

AMSC 607 / CMSC 878o Advanced Numerical Optimization

Fall 2001

UNIT 3: Constrained Optimization

PART 3: Penalty and Barrier Methods

Dianne P. O'Leary

©2001,2003

Penalty and Barrier Methods

Reference: N&S Chapter 16

Two disadvantages of feasible direction methods:

- the need to guess the active set.
- the need to get an initial feasible point.

So we again build on our unconstrained algorithms, but in a different way.

Idea:

- We have a collection of methods that work for unconstrained problems:
 - Newton's method
 - steepest descent
 - quasi-Newton
 - ...
- Can we modify them to work when we have constraints?

The plan

- Barrier Methods
- Overcoming Ill-Conditioning in Barrier Methods
- Penalty Methods
- Overcoming Ill-Conditioning in Penalty Methods: Exact Penalty Methods

Barrier Methods

Barrier Methods

Idea: Suppose we have an initial feasible point.

We want to change the function f to **raise a barrier at the boundary**.

Picture

The formalism

The problem:

$$\begin{aligned} \min_x f(x) \\ c(x) \geq 0 \end{aligned}$$

The Lagrangian:

$$L(x, \lambda) = f(x) - \lambda^T c(x)$$

Barrier function:

$$B_\mu(x) = f(x) + \mu\phi(x)$$

Possible choices of ϕ :

- **Log barrier:**

$$\phi(x) = - \sum_{i=1}^m \log(c_i(x))$$

- **Inverse barrier:**

$$\phi(x) = \sum_{i=1}^m \frac{1}{c_i(x)}$$

The severity of the barrier is controlled by the choice of μ :

- large μ : gradual barrier
- small μ : sharp barrier

Picture

The tradeoff

- **If μ is large and the minimizer is near the boundary:**
 - the barrier function looks very different from f .
- **If μ is small and the minimizer is near the boundary:**
 - steep gradients.
 - ill-conditioning and therefore hard to solve the problem.

Nice example in N&S p.534.

Because of this, we usually solve a **sequence** of problems:

- Start with a large μ to guide us near the solution.
- Gradually reduce μ , using our previous solution as a starting point.

Typical Convergence Result

Theorem: Under the conditions in N&S pp540-541 (continuity of the functions, bounded level sets, nonempty feasible set with some regularity) , let $x(\mu)$ solve the barrier problem

$$\min_x B_\mu(x).$$

Then given a sequence $\mu_1 > \mu_2 > \dots$ converging to zero, there exists a subsequence of $\{x_{\mu_i}\}$ that converges to a local solution of our problem.

Proof: See book.

Even more important:

If some technical conditions hold (a constraint qualification, and no “accidentally zero” Lagrange multiplier), then the points $x(\mu)$ define a path called the **barrier trajectory**.

And this path is differentiable!

(We will use this later when we develop **interior point methods**.)

For a proof of this result, see N&S p. 532.

A closer look at the log barrier formulation

$$B_\mu(x) = f(x) - \mu \sum_{i=1}^m \log(c_i(x))$$

Notation: To eliminate clutter in the equations, I will abbreviate the summation:

$$\sum_{i=1}^m = \sum .$$

To minimize B_μ with respect to x , we set the derivative of this **barrier function** equal to zero:

$$\begin{aligned} 0 &= g(x) - \mu \sum \frac{\nabla c_i(x)}{c_i(x)} \\ &= g(x) - \sum \frac{\mu}{c_i(x)} \nabla c_i(x) \end{aligned}$$

But in order to solve our problem, what we really want to do is set the derivative of the **Lagrangian** to zero:

$$0 = g(x) - A(x)^T \lambda = g(x) - \sum \lambda_i \nabla c_i(x)$$

So we see that we have estimates of the Lagrange multipliers:

$$\lambda_i \approx \frac{\mu}{c_i(x)} \geq 0.$$

This is good, but we haven't solved our original problem, since

$$\lambda_i(\mu)c_i(x) = \mu$$

instead of equaling zero.

So, as μ gets small, we haven't solved our original problem, but **we have solved a nearby one** formed by changing 0 to μ !

A closer look at the inverse barrier function

$$B_\mu(x) = f(x) + \mu \sum \frac{1}{c_i(x)}$$

Unquiz: Go through a similar derivation for this barrier function, showing that the Lagrange multiplier estimates are

$$\lambda_i(\mu) = \frac{\mu}{(c_i(x))^2}.$$

□

Summary of the Barrier formulation: Log barrier function

The solution $x(\mu)$ to

$$\min_x B_\mu(x) = \min_x f(x) - \mu \sum \log(c_i(x))$$

is the solution to

$$\begin{aligned} g(x) - A(x)^T \lambda &= 0 \\ \lambda &\geq 0 \\ \lambda_i c_i(x) &= \mu, \quad i = 1, \dots, m \end{aligned}$$

and this is a perturbation of order μ of the optimality conditions for our original problem.

Properties of this formulation:

- We don't need to choose an active set, so there is **no combinatorial explosion** as the number of constraints gets large.

- We have convergence under rather mild conditions.
- We get estimates of the Lagrange multipliers free.
- B is convex if f and $-c_i$ are.
- The Hessian matrix of B is ill-conditioned as $\mu \rightarrow 0$. This makes it hard to compute the Newton direction, but we'll fix this in a minute.
- The barrier function is rather ill behaved for small μ so the line search must be specially designed:
 - A quadratic model does not fit the function well.
 - We need to model the singularity in the function.
- We need a **strictly feasible** initial guess.

Overcoming Ill-Conditioning in Barrier Methods

Note: This covers material similar to N&S 16.3 but is somewhat different.

We'll use the log barrier function as an example, and consider what we would need to do in order to compute the Newton direction.

The log barrier function:

$$B_\mu(x) = f(x) - \mu \sum \log(c_i(x))$$

Differentiate with respect to x :

$$\nabla B_\mu = g(x) - \mu \sum \frac{\nabla c_i(x)}{c_i(x)}$$

Calculate the second derivative matrix, recalling that $\lambda_i = \mu/c_i(x)$:

$$\begin{aligned} \nabla^2 B_\mu &= \mathbf{H} - \mu \sum \left[\frac{\nabla^2 c_i(\mathbf{x})}{c_i(\mathbf{x})} - \frac{\nabla c_i(\mathbf{x}) \nabla c_i(\mathbf{x})^T}{c_i(\mathbf{x})^2} \right] \\ &= \mathbf{H} - \sum \lambda_i \nabla^2 c_i(\mathbf{x}) + \frac{1}{\mu} \sum \lambda_i^2 \nabla c_i(\mathbf{x}) \nabla c_i(\mathbf{x})^T \\ &\equiv \mathbf{H}_L + \frac{1}{\mu} \mathbf{A}^T \mathbf{D} \mathbf{A} \end{aligned}$$

where D is a diagonal matrix with entries λ_i^2 .

The first term, H_L , is the **Hessian of the Lagrangian function**, and this is **independent of μ** .

The second term is troublesome. If the i th constraint is active at the solution, then as $\mu \rightarrow 0$, λ_i will be nonzero, so $\lambda_i^2/\mu \rightarrow \infty$.

This causes the Hessian matrix for B_μ to blow up in $k < n$ different directions if k constraints are active at the solution. This is not good for Newton's method....

Linear algebra to the rescue

For simplicity, assume that all constraints are active at the solution, and all $\lambda_i > 0$.

Factor

$$A^T = [Q_1 \quad Q_2] \begin{bmatrix} R \\ 0 \end{bmatrix}$$

where

- $A^T : n \times m$
- $Q_1 : n \times m$
- $Q_2 : n \times (n - m)$
- $R : m \times m$
- $Q^T Q = I$.

Then

$$\nabla^2 B_\mu = H_L + \frac{1}{\mu} A^T D A = H_L + \frac{1}{\mu} Q_1 R D R^T Q_1^T.$$

Now let

$$\tilde{Q} = [Q_2 \quad Q_1].$$

Then

$$\begin{aligned} \tilde{Q}^T A^T D A \tilde{Q} &= \begin{bmatrix} Q_2^T \\ Q_1^T \end{bmatrix} Q_1 R D R^T Q_1^T [Q_2 \quad Q_1] \\ &= \begin{bmatrix} 0 & 0 \\ 0 & R D R^T \end{bmatrix}, \end{aligned}$$

so

$$\begin{aligned}\tilde{Q}^T \nabla^2 B_\mu \tilde{Q} &= \tilde{Q}^T H_L \tilde{Q} + \frac{1}{\mu} \begin{bmatrix} 0 & 0 \\ 0 & RDR^T \end{bmatrix} \\ &\equiv \tilde{H} + \frac{1}{\mu} \begin{bmatrix} 0 & 0 \\ 0 & RDR^T \end{bmatrix}\end{aligned}$$

Now let's factor this matrix as

$$\tilde{Q}^T \nabla^2 B_\mu \tilde{Q} = \begin{bmatrix} I & 0 \\ K & I \end{bmatrix} \begin{bmatrix} \tilde{H}_{11} & \tilde{H}_{12} \\ 0 & G \end{bmatrix}$$

where \tilde{H}_{ij} refers to a block of \tilde{H} .

The first block row of this is clearly ok.

Setting the (2,1) block on the left-hand side equal to the block on the right-hand side yields

$$\tilde{H}_{21} = K \tilde{H}_{11}$$

so

$$K = \tilde{H}_{21} \tilde{H}_{11}^{-1}$$

and the conditioning of \tilde{H}_{11}^{-1} is not affected by μ .

Next, look at the (2,2) block:

$$\tilde{H}_{22} + \frac{1}{\mu} RDR^T = K \tilde{H}_{12} + G$$

so

$$\begin{aligned}G &= \frac{1}{\mu} RDR^T + \tilde{H}_{22} - \tilde{H}_{21} \tilde{H}_{11}^{-1} \tilde{H}_{12} \\ &\rightarrow \frac{1}{\mu} RDR^T \\ &\equiv \hat{G}\end{aligned}$$

as $\mu \rightarrow 0$. **Note that the condition of \hat{G} is independent of μ !**

So we will replace G by \hat{G} in the factorization to get an approximate Newton direction.

Where are we?

We want to solve

$$\nabla^2 B_\mu p = y$$

where y is the negative gradient of the barrier function.

This is equivalent to solving

$$\tilde{Q}^T \nabla^2 B_\mu \tilde{Q} (\tilde{Q}^T p) = \tilde{Q}^T y.$$

How can we (**approximately**) solve this?

Algorithm:

1. Let $\tilde{y} = \tilde{Q}^T y$.

2. Solve

$$\begin{bmatrix} I & 0 \\ K & I \end{bmatrix} z = \tilde{y}.$$

Since $K\tilde{H}_{11} = \tilde{H}_{21}$, we can solve this by forming

$$\begin{aligned} z_1 &= \tilde{y}_1 \\ \tilde{H}_{11}q &= z_1 \\ z_2 &= \tilde{y}_2 - \tilde{H}_{21}q \end{aligned}$$

3. Solve

$$\begin{bmatrix} \tilde{H}_{11} & \tilde{H}_{12} \\ 0 & \hat{\mathbf{G}} \end{bmatrix} \tilde{p} = z$$

by forming

$$\begin{aligned} \hat{\mathbf{G}}\tilde{p}_2 &= z_2 \\ \tilde{H}_{11}\tilde{p}_1 &= z_1 - \tilde{H}_{12}\tilde{p}_2 \end{aligned}$$

4. Form $p = \tilde{Q}\tilde{p}$.

Then p is the **approximate** Newton direction and ill-conditioning has been avoided!

We have already developed most of the machinery we need to understand these methods, and they are not as important in practice as the Barrier methods, so we'll cover them in somewhat less detail.

Penalty Methods

If we don't have an initial feasible point, then none of our previous algorithms (feasible direction methods, barrier methods) can be applied.

What can we do? **Use a penalty method.**

Idea: Let

$$\pi(x, \rho) = f(x) + \rho\psi(x)$$

where ρ is a scalar **penalty parameter** and

$$\psi(x) = \begin{cases} 0 & x \text{ feasible} \\ > 0 & \text{otherwise} \end{cases}$$

Again we solve a sequence of problems

$$\min_x \pi(x, \rho_i),$$

but now $\rho_i \rightarrow \infty$.

Examples of penalty functions:

- **Quadratic loss function for equality constraints:**

$$\psi(x) = \frac{1}{2}c(x)^T c(x)$$

- **Quadratic loss function for inequalities:**

$$\psi(x) = \frac{1}{2}c_+(x)^T c_+(x)$$

where

$$(c_+)_i = \begin{cases} c_i & c_i < 0 \\ 0 & c_i \geq 0 \end{cases}$$

- One other loss function:

$$\psi(x) = \frac{1}{\gamma} \sum |(c_+(x))_i|^\gamma$$

for a fixed parameter $\gamma \geq 1$.

Properties

- The sequence of optimal points $\{x(\rho_i)\}$ is convergent (under mild assumptions)
- The points $x(\rho)$ define a trajectory and yield estimates of Lagrange multipliers.

Example: quadratic loss function for equality constraints

$$\pi(x, \rho) = f(x) + \frac{1}{2}\rho \sum c_i(x)^2$$

Differentiate:

$$\nabla \pi = g(x) + \rho \sum c_i(x) \nabla c_i(x) = 0$$

Compare with Lagrangian conditions:

$$0 = g(x) - \sum \lambda_i \nabla c_i(x),$$

so our estimates of Lagrange multipliers are

$$\lambda_i(\rho) = -\rho c_i(x).$$

□

Properties:

- The conditioning of the Hessian of the penalty function is increasingly bad as $\rho \rightarrow \infty$. (See the example in N&S p.538.)

- For inequality constraints, the second derivative can be discontinuous. This gives Newton's method a lot of trouble!

Overcoming Ill-Conditioning in Penalty Methods: Exact Penalty Methods

Reference: N&S 16.5.

Idea: Construct a penalty problem that is **equivalent** to the original problem.

Then we don't need to solve a sequence of problems!

Choice 1

Sacrifice differentiability.

$$\pi(x, \rho) = f(x) + \rho \sum |(c_+(x))_i|$$

This function is continuous, but fails to be differentiable when we hit a constraint.

Lemma: If x^* satisfies the 2nd order sufficient conditions for optimality, then there exists a number $\bar{\rho}$ such that x^* is an isolated local minimizer of $\pi(x, \rho)$ for any ρ greater than $\bar{\rho}$.

Idea of Proof: See picture. Choose ρ large enough so that $\pi(x, \rho) > f(x^*)$ for infeasible points in the neighborhood, and so that x^* is the local minimizer for feasible points. \square

A subtle point: We can't make this work with a quadratic penalty function – it is too flat at the constraint. See the example in N&S p538.

Choice 2

Put the Lagrange multipliers into the function explicitly: **Augmented Lagrangian**

We'll use the original function, plus a penalty term, plus the Lagrangian term:

$$\min_x f(x) + \frac{1}{2} \rho c(x)^T c(x) - \lambda^T c(x)$$

We need to choose ρ large enough to make x^* an unconstrained minimizer of this function.

The resulting algorithm:

Given initial guesses $x^{(0)}$, $\lambda^{(0)}$, and $\rho^{(0)}$, set $k = 0$.

We iterate until optimality:

- Determine $x^{(k+1)}$ by minimizing the augmented Lagrangian, fixing $\lambda = \lambda^{(k)}$ and $\rho = \rho^{(k)}$, using your favorite method.
- Get new estimates of the multipliers $\lambda^{(k+1)}$ and increase $\rho^{(k)}$ to get $\rho^{(k+1)}$.
- Set $k = k + 1$.

How to estimate the multipliers: Differentiate the augmented Lagrangian that we just minimized, and evaluate it at $x^{(k+1)}$:

$$g(x^{(k+1)}) - A(x^{(k+1)})^T (\lambda^{(k)} - \rho^{(k)} c(x^{(k+1)})) = 0,$$

so we set

$$\lambda^{(k+1)} = \lambda^{(k)} - \rho^{(k)} c(x^{(k+1)}).$$

An alternate interpretation of the augmented Lagrangian method

The iteration is really driven by λ , not ρ .

$$\lambda^{(k+1)} = \lambda^{(k)} - \rho^{(k)} c(x^{(k+1)})$$

means that we change λ by using the **search direction** $c(x^{(k+1)})$ and the **step length parameter** $\rho^{(k)}$.

This is **steepest ascent** on the dual problem

$$\max_{\lambda} \min_x f(x) - \lambda^T c(x) + \frac{1}{2} \rho c(x)^T c(x)$$

since if we differentiate $f(x) - \lambda^T c(x) + \frac{1}{2} \rho c(x)^T c(x)$ respect to λ we get

$$-c(x(\lambda)).$$

For inequality constraints:

$$\min_x f(x) - \mu^{(k)} \sum (\lambda^{(k)})_i \log(\mu^{(k)-1} c_i(x) + 1)$$

(See N&S p560)

How do we update λ ?

$$g(x^{(k+1)}) - \mu^{(k)} \sum (\lambda^{(k)})_i \frac{1}{\mu^{(k)-1} c_i(x^{(k+1)}) + 1} \mu^{(k)-1} \nabla c_i(x^{(k+1)}) = 0$$

so

$$(\lambda^{(k+1)})_i = \frac{(\lambda^{(k)})_i}{\mu^{(k)-1} c_i(x^{(k+1)}) + 1}.$$

For augmented Lagrangians ...

... there is a tension between conditioning and convergence:

- fast if $\rho^{(k)}$ large ,
- but ill-conditioned, so hard to get a good solution.

Final Words

These methods still have their place, but their main usefulness to us is as forerunners of **interior point methods**.