# Algorithms for Constrained Optimization

## ME598/494 Lecture

## Max Yi Ren

Department of Mechanical Engineering, Arizona State University

April 5, 2016

# Outline

1. Convergence
2. Sequential quadratic programming (SQP)
3. Quasi-Newton Methods
4. Active set strategies
5. Penalty and barriers
6. Augmented Lagrangian

# Global and local convergence

**Global convergence** refers to the ability of the algorithm to reach the neighborhood of some local solution $\mathbf{x}_*$ from an arbitrary initial point $\mathbf{x}_0$, which is not close to $\mathbf{x}_*$. The convergence of a globally convergent algorithm should not be affected by the choice of initial point.

**Local convergence** refers to the ability of the algorithm to approach $\mathbf{x}_*$, rapidly from a point in the neighborhood of $\mathbf{x}_*$.

**convergence ratio** $\gamma$:

- Linear convergence: $||\mathbf{x}_{k+1} - \mathbf{x}_*|| \leq \gamma ||\mathbf{x}_k - \mathbf{x}_*||, 0 < \gamma < 1$
- Quadratic convergence: $||\mathbf{x}_{k+1} - \mathbf{x}_*|| \leq \gamma ||\mathbf{x}_k - \mathbf{x}_*||^2, \gamma \in \mathbb{R}$

Newton's method has quadratic convergence rate but is not globally convergent; Gradient descent has global convergence but in some cases can be inefficient.

# The Lagrange-Newton Equations (1/2)

Consider the equality constrained problem

$$\begin{aligned}\min \quad & f(\mathbf{x}) \\ \text{subject to} \quad & \mathbf{h}(\mathbf{x}) = \mathbf{0}\end{aligned} \tag{1}$$

The stationary condition for this problem is

$$\nabla L(\mathbf{x}_*, \boldsymbol{\lambda}_*) = \mathbf{0}^T.$$

We may solve this equation using Newton-Ralphson to update $\mathbf{x}$ and $\boldsymbol{\lambda}_*$:

$$[\nabla L(\mathbf{x}_k + \partial\mathbf{x}_k, \boldsymbol{\lambda}_k + \partial\boldsymbol{\lambda}_k)]^T = \nabla L_k^T + \nabla^2 L_k [\partial\mathbf{x}_k, \partial\boldsymbol{\lambda}_k]^T, \tag{2}$$

where

$$\nabla L_k^T = [\nabla f_k^T + \nabla \mathbf{h}_k^T \boldsymbol{\lambda}_k, \ \mathbf{h}^T], \tag{3}$$

and

$$\nabla^2 L_k = [\nabla^2 f + \boldsymbol{\lambda}^T \nabla^2 \mathbf{h}, \ \nabla \mathbf{h}^T; \ \nabla \mathbf{h}, \ \mathbf{0}]_k \tag{4}$$

.

## The Lagrange-Newton Equations (2/2)

Define $\mathbf{W} = \nabla^2 f + \boldsymbol{\lambda}^T \nabla^2 \mathbf{h}$ and $\mathbf{A} = \nabla \mathbf{h}$ to have

$$\nabla^2 L_k = [\mathbf{W} \ \mathbf{A}^T; \mathbf{A} \ \mathbf{0}]_k.$$

Denote the step as $\mathbf{s}_k := \partial \mathbf{x}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ and set the left-hand side of Equation (2) to zero to have

$$\begin{aligned}
\mathbf{W}_k \mathbf{s}_k + \mathbf{A}_k^T \boldsymbol{\lambda}_{k+1} + \nabla f_k^T &= \mathbf{0} \\
\mathbf{A}_k \mathbf{s}_k + \mathbf{h}_k &= \mathbf{0}
\end{aligned} \tag{5}$$

Equation (5) is referred to as a Lagrange-Newton method for solving the constrained problem (1).

What are the conditions of $\mathbf{W}_*$ and $\mathbf{A}_*$ for the solution to be unique?

## Quadratic programming subproblem

Note that Equation (5) can be viewed as the KKT conditions for the quadratic model

$$
\begin{aligned}
\min_{\mathbf{s}_k} \quad & q(\mathbf{s}_k) = f_k + \nabla_{\mathbf{x}} L_k \mathbf{s}_k + \frac{1}{2} \mathbf{s}_k^T \mathbf{W}_k \mathbf{s}_k \\
\text{subject to} \quad & \mathbf{A}_k \mathbf{s}_k + \mathbf{h}_k = \mathbf{0},
\end{aligned}
\tag{6}
$$

where $\nabla_{\mathbf{x}} L_k = \nabla f_k + \boldsymbol{\lambda}_k^T \nabla \mathbf{h}_k$ and the multiplier of problem (6) are $\partial \boldsymbol{\lambda}_k$.

It can be shown that solving the Lagrange-Newton equations from Equation (5) is equivalent to solving the quadratic programming subproblem (6).

An alternative QP subproblem is as follows

$$
\begin{aligned}
\min_{\mathbf{s}_k} \quad & q(\mathbf{s}_k) = f_k + \nabla f_k \mathbf{s}_k + \frac{1}{2} \mathbf{s}_k^T \mathbf{W}_k \mathbf{s}_k \\
\text{subject to} \quad & \mathbf{A}_k \mathbf{s}_k + \mathbf{h}_k = \mathbf{0},
\end{aligned}
\tag{7}
$$

which also gives a solution $\mathbf{s}_k$ with multipliers $\boldsymbol{\lambda}_{k+1}$ directly, rather than $\partial \boldsymbol{\lambda}_k$. What is the meaning of this QP subproblem?

# SQP Algorithm (without line search)

1. Select initial point $\mathbf{x}_0$, $\boldsymbol{\lambda}_0$; let $k = 0$.
2. Solve the QP subproblem and determine $\mathbf{s}_k$ and $\boldsymbol{\lambda}_{k+1}$.
3. Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$. Set $k = k + 1$.
4. If KKT condition not satisfied, return to 2.

Advantage:

▶ Simple

▶ Fast, locally quadratic convergence

## Enhancements of the basic algorithm

The basic SQP algorithm may not have global convergence. For points far from $\mathbf{x}_*$, the QP subproblem may have an unbounded solution.

It is shown that for the QP subproblem to have a well-defined solution, the following is needed:

- $\mathbf{A}$ has full rank
- $\mathbf{W}$ has to be positive definite in feasible perturbations

One possibility is to use the QP solution $\mathbf{s}_k$ as a search direction and find the step size $\alpha_k$ that minimizes a *merit function*, which is a penalty function that properly weighs objective function decrease and constraint violations.

## Line search on the merit

One merit function (exact penalty function, Powell 1978a) that is widely implemented has the following form

$$\phi(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^{m_1} w_j |h_j| + \sum_{j'=1}^{m_2} \bar{w}_{j'} \max\{0, g_j\},$$

where $m_1$ and $m_2$ are the numbers of equality and inequality constraints and $w_j$ ($\bar{w}$) are weights used to balance the infeasibilities. The suggested values are

$w_{j,k} = |\lambda_{j,k}|$   for $k = $ 0th SQP iteration, and $j$th equality constraint,

$w_{j,k} = \max\{|\lambda_{j,k}|, 0.5(w_{j,k-1} + |\lambda_{j,k}|)\}$   for $k \geq 1$.

For $\bar{w}_{j',k}$, we use the same update, with $\mu_j$ in place of $\lambda_j$.

One can also use a quadratic penalty function

$$\phi(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^{m_1} w_j h_j^2 + \sum_{j=1}^{m_2} w_j (\max\{0, g_j\})^2.$$

# SQP Algorithm (with line search)

1. Select initial point $\mathbf{x}_0$, $\boldsymbol{\lambda}_0$; let $k = 0$.

2. Solve the QP subproblem and determine $\mathbf{s}_k$ and $\boldsymbol{\lambda}_{k+1}$ (and $\boldsymbol{\mu}_{k+1}$).

3. Line search: Input $\mathbf{x}_k$, $\boldsymbol{\lambda}_k$, $w_{j,k-1}$ (and $\bar{w}_{j',k-1}$); Output $w_{j,k}$, $\alpha_k$

   3.1 Calculate $w_{j,k}$ (and $\bar{w}_{j',k}$) based on $\lambda_{j,k}$ and $w_{j,k-1}$ (and $\bar{w}_{j',k-1}$)

   3.2 Calculate Jacobian matrices for $h$ and $g$, at $\mathbf{x}_k$

   3.3 Do Amijo line search and get $\alpha_k$

4. Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k$. Set $k = k + 1$.

5. If KKT condition not satisfied, return to 2.

# Quasi-Newton Methods

1. DFP
2. BFGS

## Quasi-Newton methods

Recall the Newton's method

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{H}_k^{-1} \mathbf{g}_k,$$

where $\mathbf{H}_k$ and $\mathbf{g}_k$ are the Hessian and gradient at iteration $k$. The method has quadratic convergence when $\mathbf{H}_k$ remains positive-definite. Quasi-Newton methods build a positive-definite approximation of the Hessian using $f_k$ and $\mathbf{g}_k$, and is regarded as the best general-purpose methods for solving unconstrained problems.

Calculating $\mathbf{H}_k$ can be time consuming. Therefore we wish to approximate $\mathbf{H}_k$ as $\hat{\mathbf{H}}_k$ iteratively:

$$\hat{\mathbf{H}}_{k+1} = \hat{\mathbf{H}}_k + \text{something,}$$

to get the second-order approximation at $\mathbf{x}_{k+1}$:

$$f(\mathbf{x}) = f(\mathbf{x}_{k+1}) + \mathbf{g}_{k+1}^T \partial \mathbf{x}_{k+1} + \frac{1}{2} \partial \mathbf{x}_{k+1}^T \hat{\mathbf{H}}_{k+1} \partial \mathbf{x}_{k+1}, \tag{8}$$

where $\partial \mathbf{x}_{k+1} = \mathbf{x} - \mathbf{x}_{k+1}$.

# The DFP method (1/3)

Three conditions need to be imposed on Equation (8).

First, $\hat{\mathbf{H}}$ needs to be symmetric and positive-definite.

Second, the approximated $f(\mathbf{x})$ must match the true gradients at $\mathbf{x}_k$ and $\mathbf{x}_{k+1}$.
For $\mathbf{x}_{k+1}$, the approximation from Equation (8) naturally follows that

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}|_{k+1} = \mathbf{g}_{k+1}^T.$$

Therefore the approximated gradient is the true gradient at $\mathbf{x}_{k+1}$.

For $\mathbf{x}_k$, considering a general search $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k$, we have

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}|_k = \mathbf{g}_{k+1}^T - \alpha_k \mathbf{s}_k^T \mathbf{H}_{k+1}.$$

By rearranging terms we have

$$\mathbf{g}_{k+1} - \mathbf{g}_k = \alpha_k \mathbf{H}_{k+1} \mathbf{s}_k \qquad (9)$$

Equation (9) is called the **secant equation** and key to approximate $\mathbf{H}_{k+1}$.

# The DFP method (2/3)

Multiply both ends of the secant equation by $\mathbf{s}_k^T$ to have

$$\mathbf{s}_k^T(\mathbf{g}_{k+1} - \mathbf{g}_k) = \alpha_k \mathbf{s}_k^T \mathbf{H}_{k+1} \mathbf{s}_k > 0,$$

when the Hessian is positive-definite. Note that $\mathbf{s}_k^T(\mathbf{g}_{k+1} - \mathbf{g}_k) > 0$ is called the curvature condition.

In fact, the curvature condition ensures a positive-definite approximation $\hat{\mathbf{H}}$ of $\mathbf{H}$.

However, there are infinitely many symmetric positive-definite matrices that satisfy the secant equation.

The last condition: We will select $\hat{\mathbf{H}}_{k+1}$ that is closest to $\hat{\mathbf{H}}_k$ in the weighted Frobenius norm. Overall, we find $\hat{\mathbf{H}}_{k+1}$ that solves the following convex problem

$$\min_{\mathbf{H}} \quad ||\hat{\mathbf{H}} - \hat{\mathbf{H}}_k||_F$$
$$\text{subject to} \quad \hat{\mathbf{H}} = \hat{\mathbf{H}}^T$$
$$\mathbf{g}_{k+1} - \mathbf{g}_k = \alpha_k \hat{\mathbf{H}} \mathbf{s}_k$$

# The DFP method (3/3)

Solve Problem (14) and denote $\mathbf{B} = \mathbf{H}^{-1}$ to have the DFP update

$$\mathbf{B}_{k+1}^{\text{DFP}} = \mathbf{B}_k + \left[\frac{\partial\mathbf{x}\partial\mathbf{x}^T}{\partial\mathbf{x}^T\partial\mathbf{g}}\right]_k - \left[\frac{(\mathbf{B}\partial\mathbf{g})(\mathbf{B}\partial\mathbf{g})^T}{\partial\mathbf{g}^T\mathbf{B}\partial\mathbf{g}}\right]_k, \tag{10}$$

where $\partial\mathbf{x}_k = \alpha\mathbf{s}_k$ and $\partial\mathbf{g} = \mathbf{g}_{k+1} - \mathbf{g}_k$.

The Davidon-Fletcher-Powell (DFP) method was originally proposed by W.C. Davidon in 1959.

# The BFGS method (1/2)

Instead of imposing conditions on the Hessian as in DFP, the BFGS method directly work on the inverse of the Hessian. The secant equation therefore is in the form

$$\mathbf{B}_{k+1}(\mathbf{g}_{k+1} - \mathbf{g}_k) = \alpha_k \mathbf{s}_k.$$

The revised conditions lead to the BFGS update

$$\mathbf{B}_{k+1}^{BFGS} = \mathbf{B}_k + \left[1 + \frac{\partial \mathbf{g}^T \mathbf{B} \partial \mathbf{g}}{\partial \mathbf{x}^T \partial \mathbf{g}}\right]_k \left[\frac{\partial \mathbf{x} \partial \mathbf{x}^T}{\partial \mathbf{x}^T \partial \mathbf{g}}\right]_k - \left[\frac{\partial \mathbf{x} \partial \mathbf{g}^T \mathbf{B} + \mathbf{B} \partial \mathbf{g} \partial \mathbf{x}^T}{\partial \mathbf{x}^T \partial \mathbf{g}}\right]_k.$$

The Broyden-Fletcher-Goldfarb-Shanno (BFGS) method is more commonly used than DFP.

# The BFGS method (2/2)

There are several ways to set the initial value for $\mathbf{H}$ (or $\mathbf{B}$):

- A finite difference approximation at $\mathbf{x}_0$.
- Use the identity matrix.
- Use $\text{diag}(\lambda_1, \lambda_2, ...)$, where $\boldsymbol{\lambda}$ captures the scaling of the variables.

# BFGS for SQP

The evaluation of $\mathbf{W}_k$ can be approximated using BFGS. Note that the curvature condition is not necessarily satisfied by line search. The following adjustment keeps Hessian approximation positive-definite, i.e., $\partial\mathbf{x}_k^T\partial\mathbf{g}_k > 0$.

$$\partial\mathbf{g}_k = \theta_k\mathbf{y}_k + (1 - \theta_k)\hat{\mathbf{W}}_k\partial\mathbf{x}_k, \quad 0 \leq \theta \leq 1,$$

where

$$\mathbf{y}_k = \nabla L(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1})^T - \nabla L(\mathbf{x}_k, \boldsymbol{\lambda}_{k+1})^T,$$

and

$$\theta_k = \begin{cases} 1 & \text{if } \partial\mathbf{x}_k^T\mathbf{y}_k \geq (0.2)\partial\mathbf{x}_k^T\hat{\mathbf{W}}_k\partial\mathbf{x}_k, \\ \frac{(0.8)\partial\mathbf{x}_k^T\hat{\mathbf{W}}_k\partial\mathbf{x}_k}{\partial\mathbf{x}_k^T\hat{\mathbf{W}}_k\partial\mathbf{x}_k - \partial\mathbf{x}_k^T\mathbf{y}_k} & \text{if } \partial\mathbf{x}_k^T\mathbf{y}_k < (0.2)\partial\mathbf{x}_k^T\hat{\mathbf{W}}_k\partial\mathbf{x}_k, \end{cases}$$

and $\partial\mathbf{x}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$, $\hat{\mathbf{W}}_k$ is the current BFGS approximation to the Hessian of the Lagrangian.

# SQP Algorithm (with line search, BFGS)

1. Select initial point $\mathbf{x}_0$, $\boldsymbol{\lambda}_0$, $\mathbf{W}$; let $k = 0$.

2. Solve the QP subproblem and determine $\mathbf{s}_k$ and $\boldsymbol{\lambda}_{k+1}$ (and $\boldsymbol{\mu}_{k+1}$).

3. Line search: Input $\mathbf{x}_k$, $\boldsymbol{\lambda}_k$, $w_{j,k-1}$ (and $\bar{w}_{j',k-1}$); Output $w_{j,k}$, $\alpha_k$

    3.1 Calculate $w_{j,k}$ (and $\bar{w}_{j',k}$) based on $\lambda_{j,k}$ and $w_{j,k-1}$ (and $\bar{w}_{j',k-1}$)

    3.2 Calculate Jacobian matrices for $h$ and $g$, at $\mathbf{x}_k$

    3.3 Do Amijo line search and get $\alpha_k$

4. Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k$.

5. BFGS: Input $\mathbf{W}_k$, $\mathbf{x}_k$, $\mathbf{x}_{k+1}$, $\boldsymbol{\lambda}_{k+1}$; Output $\mathbf{W}_{k+1}$

    5.1 Calculate $\mathbf{y}_k$, $\theta_k$ and $\mathbf{g}_k$

    5.2 Calculate $\mathbf{W}_{k+1}$

6. Set $k = k + 1$.

7. If KKT condition not satisfied, return to 2.

## Active set strategy

We now discuss how to deal with *inequality* constraints in SQP (and GRG). The difficulty is that we do not know at the beginning which inequality constraints will be active at an optimal solution. The strategy is to maintain a *working set* of active constraints (along with equality constraints) and keep adding or deleting constraints in this working set.
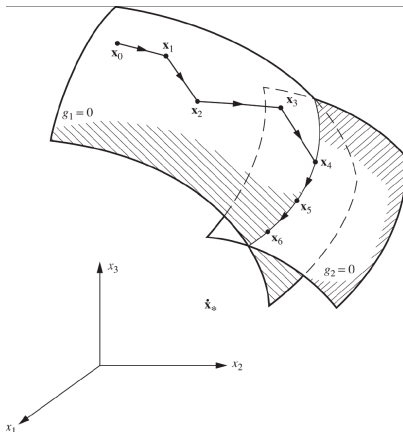
## Active set for SQP

Active set can be applied in two ways

- ▶ An active set strategy may be employed on the original problem so that the QP subproblem always have only equality constraints.

- ▶ The second way is to pose the QP subproblem with the linearized inequalities included ($\mathbf{A}_k \mathbf{s}_k + \mathbf{g}_k \leq \mathbf{0}$), and use an active set strategy on the subproblem.

We discuss the second approach.

Note: The merit function must then include all constraints, active and inactive, to guard against failure when the wrong active set is used to determine the step direction.

## Adding constraints

Starting at an initial feasible point (e.g., a zero vector) and an initial working set (e.g., only equality constraints), we solve the QP subproblem subject to the linearized equalities in the working set. When hitting a new inequality constraint (moving from $\mathbf{x}_3$ to $\mathbf{x}_4$ in the figure), that constraint will be added to the working set and the step size is reduced to retain feasibility.

## Removing constraints

When arrived at a point where no progress is possible by adding constraints, we check the KKT conditions and estimate the Lagrangian multipliers (since we may not have arrived at an optimal solution yet, these multipliers are only estimated). If the Lagrangian multipliers for some active constraints are negative, these constraints will become candidate for deletion. A common heuristic is to delete one constraint with the most negative multiplier. (Why?)

## Solving the quadratic subproblem

consider the QP problem

$$
\begin{aligned}
\min \quad & \frac{1}{2}\mathbf{s}^T\mathbf{W}\mathbf{s} + \mathbf{c}^T\mathbf{s} \\
\text{subject to} \quad & \mathbf{A}\mathbf{s} - \mathbf{b} = \mathbf{0}.
\end{aligned}
\tag{11}
$$

The Lagrange-Newton equations (KKT conditions) for this problem is

$$
\begin{pmatrix} \mathbf{W} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{s} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} -\mathbf{c} \\ \mathbf{b} \end{pmatrix}
$$

When the Lagrangian matrix is invertible (e.g., $\mathbf{W}$ positive-definite and $\mathbf{A}$ full rank), the solution to the QP problem is

$$
\begin{pmatrix} \mathbf{s} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} \mathbf{W} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} -\mathbf{c} \\ \mathbf{b} \end{pmatrix}
$$

# Active set strategy with GRG

1. Input initial feasible point and working set.

2. Compute a feasible search vector $\mathbf{s}_k$ in the reduced space.

3. Compute a step length $\alpha_k$ along $\mathbf{s}_k$, such that $f(\mathbf{x}_k + \alpha_k \mathbf{s}_k) < f(\mathbf{x}_k)$. If $\alpha_k$ violates a constraint, continue; otherwise go to 6.

4. Add a violated constraint to the constraint set and reduce $\alpha_k$ to the maximum possible value that retains feasibility.

5. Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k$.

6. Check the norm of reduced gradient. If not zero, go to step 2. Otherwise, check if estimates of Lagrangian multipliers for active constraints are positive or not. If not all positive, delete a constraint that has the most negative multiplier, and go to step 2. Otherwise, terminate.

# Active set strategy for QP subproblem in SQP

1. Input QP subproblem and working set (linearized equality constraints).

2. Solve the quadratic subproblem to get $\mathbf{s}_k$ and $\boldsymbol{\lambda}_{k+1}$ (and $\boldsymbol{\mu}_{k+1}$).

3. If inequality constraints are violated at $\mathbf{s}_k$, add the most violated one to the working set.

4. If some $\boldsymbol{\mu}_{k+1}$ are negative or zero, find the most negative one and remove it from the working set.

5. If all constraints are satisfied and all $\boldsymbol{\mu}_{k+1}$ (in the working set) are positive, terminate. Otherwise go to 2.

# SQP Algorithm (with line search, BFGS, active set)

1. Select initial point $\mathbf{x}_0$, $\boldsymbol{\lambda}_0$, $\mathbf{W}$; let $k = 0$.

2. Solve the QP subproblem: Input: QP subproblem at $\mathbf{x}_k$; Output: $\mathbf{s}_k$ and $\boldsymbol{\lambda}_{k+1}$ (and $\boldsymbol{\mu}_{k+1}$).

3. Line search: Input $\mathbf{x}_k$, $\boldsymbol{\lambda}_k$, $w_{j,k-1}$ (and $\bar{w}_{j',k-1}$); Output $w_{j,k}$, $\alpha_k$

   3.1 Calculate $w_{j,k}$ (and $\bar{w}_{j',k}$) based on $\lambda_{j,k}$ and $w_{j,k-1}$ (and $\bar{w}_{j',k-1}$)

   3.2 Calculate Jacobian matrices for $h$ and $g$, at $\mathbf{x}_k$

   3.3 Do Amijo line search and get $\alpha_k$

4. Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k$.

5. BFGS: Input $\mathbf{W}_k$, $\mathbf{x}_k$, $\mathbf{x}_{k+1}$, $\boldsymbol{\lambda}_{k+1}$; Output $\mathbf{W}_{k+1}$

   5.1 Calculate $\mathbf{y}_k$, $\theta_k$ and $\mathbf{g}_k$

   5.2 Calculate $\mathbf{W}_{k+1}$

6. Set $k = k + 1$.

7. If KKT condition not satisfied, return to 2.

# Barrier method (1/2)

Instead of solving the constrained problem, we can construct a *barrier function* to be optimized

$$T(\mathbf{x}, r) := f(\mathbf{x}) + rB(\mathbf{x}), \quad r > 0,$$

where $B(\mathbf{x}) := -\sum_{j=1}^{m} ln[-g_j(\mathbf{x})]$ (logarithmic) or $B(\mathbf{x}) := -\sum_{j=1}^{m} g_j^{-1}(\mathbf{x})$ (inverse). The barrier method only works for problems with inequality constraints.

## Barrier function algorithm

1. Find an interior point $\mathbf{x}_0$. Select a monotonically decreasing sequence $\{r_k\} \to 0$ for $k \to \infty$. Set $k = 0$.

2. At iteration $k$, minimize the function $T(\mathbf{x}, r_k)$ using an unconstrained method and $\mathbf{x}_k$ as the starting point. the solution $\mathbf{x}_*(r_k)$ is set equal to $\mathbf{x}_{k+1}$.

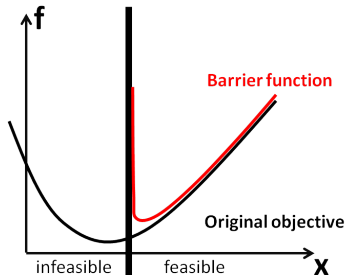3. Perform a convergence test. If the test is not satisfied, set $k = k + 1$ and return to step 2.

# Barrier method (2/2)

When using the barrier method, we can estimate Lagrange multipliers as

- $\mu_j(r_k) = -r_k/g_j$, (for the logarithmic barrier)
- $\mu_j(r_k) = -r_k/g_j^2$ (for the inverse barrier)

The actual multiplier can be obtained at the limit.

Note that this basic barrier method has a major computational difficulty: A small $r_k$ leads to an ill-conditioned Hessian, making the optimization difficult. (What can we do?)

# Penalty method (1/2)

A typical penalty function has the form

$$T(\mathbf{x}, r) := f(\mathbf{x}) + r^{-1} P(\mathbf{x}), \quad r > 0,$$

where the *penalty function* $P(\mathbf{x})$ can take a quadratic form

$$P(\mathbf{x}) := \sum_{j=1}^{m} [\max\{0, g_j(\mathbf{x})\}]^2$$

for inequality constraints, and

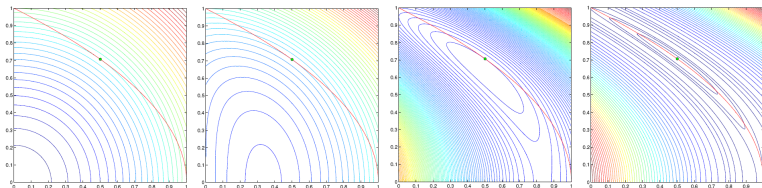$$P(\mathbf{x}) := \sum_{j=1}^{m} [h_j(\mathbf{x})]^2$$

for equality constraints.

Lagrange multipliers can be estimated as

$$\mu_j(r_k) = (2/r_k) \max\{0, g_j(\mathbf{x})_{\mathbf{k}}\}$$

for a decreasing sequence $\{r_k\}$.

# Penalty method (2/2)

An example of an ill-conditioned Hessian when using the penalty method (from lecture notes of Nick Gould)



(a) $r = 100$          (b) $r = 1$          (c) $r = 0.1$          (d) $r = 0.01$

Figure: Quadratic penalty function for min $x_1^2 + x_2^2$ subject to $x_1 + x_2^2 = 1$

# Augmented Lagrangian (1/6)

Recall that the drawback of the penalty method (as well as the barrier method) is that we can only find a good approximation of the true solution when the penalty is high, i.e., $r \to 0$, in which case the convergence of the problem will suffer from ill-conditioned Hessian.

With that in mind, we introduce the augmented Lagrangian function:

$$\Phi(\mathbf{x}, \boldsymbol{\lambda}, r) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) + \frac{1}{r}||\mathbf{h}(\mathbf{x})||^2$$

When the Lagrangian multipliers $\boldsymbol{\lambda}$ are close to their true values, a reasonable small value of $r$ allows us to find the true optimal solution $\mathbf{x}$ without encountering an ill-conditioned Hessian.

# Augmented Lagrangian (2/6)

An example where we can find the true optimal solution for a constrained problem without setting $r \to 0$. (When we guessed correctly on $\lambda_*$, we can find the solution $\mathbf{x}_*$ without $r \to 0$)
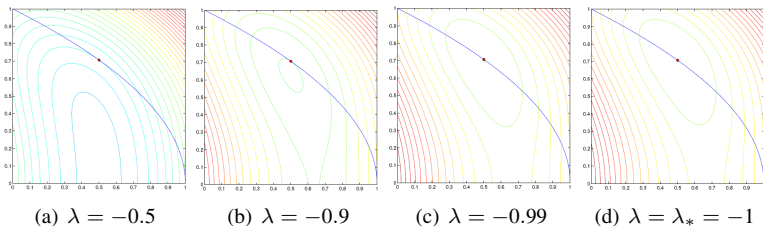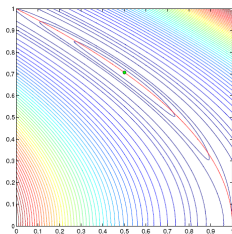


(a) $\lambda = -0.5$    (b) $\lambda = -0.9$    (c) $\lambda = -0.99$    (d) $\lambda = \lambda_* = -1$

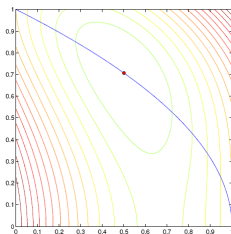Figure: Augmented Lagrangian function for min $x_1^2 + x_2^2$ subject to $x_1 + x_2^2 = 1$ with fixed $r = 1$

# Augmented Lagrangian (3/6)

Two ways of understanding augmented Lagrangian under equality constraints only:

1. **Shifted quadratic penalty function**: The augmentation shifts the origin of the penalty term so that an optimum value for the transformed function can be found without the penalty parameter going to the limit.
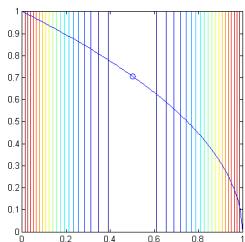


(a) $r = 0.01$        (b) $\lambda = \lambda_* = -1, r = 1$

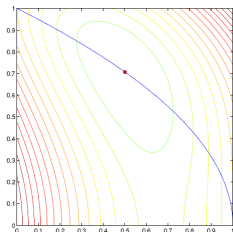Figure: Comparison between penalty and Augmented Lagrangian

# Augmented Lagrangian (4/6)

Two ways of understanding augmented Lagrangian under equality constraints only:

2. **Convexification of the Lagrangian function**: If $\mathbf{x}_*$ is a local solution for the original problem, it follows KKT conditions and the second-order sufficiency condition (which is what?). We want to construct a unconstrained problem which has $\mathbf{x}_*$ as its local solution. The first-order necessary condition of this problem should be the original KKT condition and its Hessian should be positive-definite. The augmented Lagrangian function satisfies these requirements.



(a) $\lambda = \lambda_* = -1, r = \infty$        (b) $\lambda = \lambda_* = -1, r = 1$

# Augmented Lagrangian (5/6)

The augmented Lagrangian method requires tuning of $\boldsymbol{\lambda}$ and $r$ together in some way so that $\{\mathbf{x}_k\} \to \mathbf{x}_*$.

- Check if $||\mathbf{h}(\mathbf{x})|| \leq \eta_k$ where $\{\eta_k\} \to 0$.

    - if so, set $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + 2\mathbf{h}(\mathbf{x}_k)/r_k$ and $r_{k+1} = r_k$. It is proved that such a series $\{\boldsymbol{\lambda}_k\}$ converges to $\boldsymbol{\lambda}_*$.

    - if not, set $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k$ and $r_{k+1} = \tau r_k$ for some $\tau \in (0, 1)$. Often choose $\tau = \min\{0.1, \sqrt{r_k}\}$.

- update on $\eta_k$: $\eta_k = r_k^{0.1+0.9j}$ where $j$ iterations since $r_k$ last changed.

# Augmented Lagrangian (6/6)

**The augmented Lagrangian algorithm** (from Nick Gould's lecture notes)

1. Given $r_0 > 0$ and $\boldsymbol{\lambda}_0$, set $k = 0$

2. While KKT conditions are not met

   2.1 Starting from $\mathbf{x}_k^s = \mathbf{x}_{k-1}$, use an unconstrained minimization algorithm to find an "approximate" minimizer $\mathbf{x}_k$ so that $||\nabla_{\mathbf{x}}\Phi(\mathbf{x}_k, \boldsymbol{\lambda}_k, r_k)|| \leq \varepsilon_k$

   2.2 If $||\mathbf{h}(\mathbf{x}_k)|| \leq \eta_k$, set $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + 2\mathbf{h}(\mathbf{x}_k)/r_k$ and $r_{k+1} = r_k$

   2.3 Otherwise set $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k$ and $r_{k+1} = \tau r_k$

   2.4 $k = k + 1$. Set $\varepsilon_k = r_k^{j+1}$ and $\eta_k = r_k^{0.1+0.9j}$ where $j$ iterations since $r_k$ last changed

This method can be extended to inequalities with the aid of an active set strategy. Details of implementation can be found in Pierre and Lowe (1975) and Bertsekas (1982). An alternative way proposed by Nocedal and Wright (2006) is to convert inequalities to equalities by introducing slack variables, which can be optimized separately and eliminated.

## Exercise 7.27

Using the penalty transformation $T = f + \frac{1}{2}r\mathbf{h}^T\mathbf{h}$, evaluate and sketch the progress of the penalty method (sequential unconstrained minimization) for the problem $\{\min f = x, \text{ subject to } h = x - 1 = 0\}$, with $r = 1, 10, 100, 1000$. Repeat, using the augmented Lagrangian transformation $T = f + \boldsymbol{\lambda}^T\mathbf{h} + \frac{1}{2}r\mathbf{h}^T\mathbf{h}$. (From Fletcher 1981.)

## Exercise 7.27 Solution

Penalty method: $\min x + 1/2r(x-1)^2$
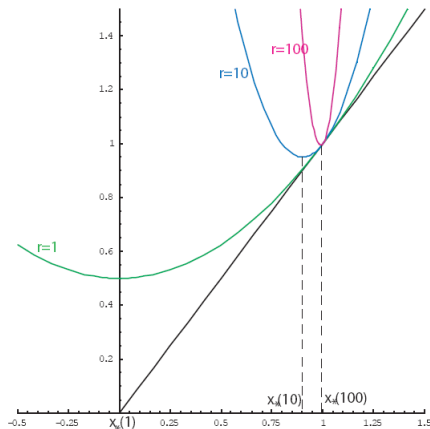
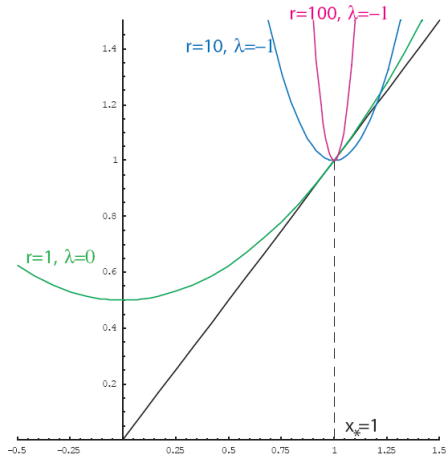Stationary point: $x_*(r) = 1 - 1/r$



Figure: Penalty method

## Exercise 7.27 Solution

Augmented Lagrangian method: $\min x + \lambda(x - 1) + 0.5r(x - 1)^2$

Stationary point: $x_*(r) = 1 - 1/r(1 + \lambda)$. Let $\lambda_0 = 0$, $r_0 = 1$, we have
$x_* = 0$, $h_0 = -1$. Then update $\lambda_1 = \lambda_0 + r_0 h_0 = -1$ and $r_1 = r_0 = 1$. Then
we have $x_* = 1$ and $h_1 = 0$.

# What have we learned so far?

- ▶ Unconstrained optimization
    - ▶ gradient descent with line search
    - ▶ Newton's method with line search
    - ▶ trust region (why?)
    - ▶ quasi-Newton (why?)
- ▶ Constrained optimization
    - ▶ generalized reduced gradient
    - ▶ Sequential quadratic programming
    - ▶ barrier and penalty (why not?)
    - ▶ augmented Lagrangian

# Comparisons

- ▶ Active set method should be attached to other algorithms and thus will not be compared with.

- ▶ GRG is the most reliable but requires the most implementation effort. It is also not the most efficient and requires a lot of function evaluation.

- ▶ Augmented Lagrangian is less reliable than GRG. A widely used package of this method is LANCELOT, which deals with large-scale optimization problems with bounds and equality constraints. The idea of augmented Lagrangian is also used in SQP type of algorithm to improve line search and Hessian approximation.

- ▶ SQP is the most widely used algorithm and can deal with large-scale problems (up to the scale of 10000 variables and constraints). It is more reliable than augmented Lagrangian and more efficient than GRG.