# Chapter 23 Algorithms for Constrained Optimization

An Introduction to Optimization
Spring, 2015

Wei-Ta Chu

# Projections

▸ The optimization algorithms described in Part II have the general form $x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$, where $d^{(k)}$ is typically a function of $\nabla f(x^{(k)})$. The value of $x^{(k)}$ is not constrained to lie inside any particular set.

▸ Consider the optimization problem
$$minimize\ f(x)$$
$$subject\ to\ \ x \in \Omega$$

▸ Modify the algorithm we already know to take constraints into account.

# Projections

▸ If $x^{(k)} + \alpha_k d^{(k)}$ is in $\Omega$, then we set form $x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$ as usual.

▸ If, on the other hand, $x^{(k)} + \alpha_k d^{(k)}$ is not in $\Omega$, then we "project" it back into $\Omega$ before setting $x^{(k+1)}$.

▸ Consider the case where the constraint set $\Omega \in R^n$ is given by

$$\Omega = \{x : l_i \leq x_i \leq u_i, \qquad i = 1, \dots, n\}$$

▸ In this case, $\Omega$ is a "box" in $R^n$; for this reason, this form of $\Omega$ is called a *box constraint.*

# Projections

▸ Given a point $\mathbf{x} \in R^n$, define $\mathbf{y} = \prod[\mathbf{x}] \in R^n$ by

$$y_i = \min\{u_i, \max\{l_i, x_i\}\} = \begin{cases} u_i & if \ x_i > u_i \\ x_i & if \ l_i \le x_i \le u_i \\ l_i & if \ x_i < l_i \end{cases}$$

▸ The point $\prod[\mathbf{x}]$ is called the projection of $\mathbf{x}$ onto $\Omega$. Note that $\prod[\mathbf{x}]$ is actually the "closest" point in $\Omega$ to $\mathbf{x}$. Using the projection operator $\prod$, we can modify the previous unconstrained algorithm as follows:

$$\mathbf{x}^{(k+1)} = \prod[\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}]$$

▸ Note that the iterates $\mathbf{x}^{(k)}$ now all lie inside $\Omega$. We call the algorithm above a *projection algorithm.*

# Projections

- In the more general case, we can define the projection onto $\Omega$:

$$\prod[x] = \arg\min_{z \in \Omega} \|z - x\|$$

- In this case, $\prod[x]$ is again the "closest" point in $\Omega$ to $x$. This projection operator is well-defined only for certain types of constraint sets: for example, closed convex sets.

- For some sets $\Omega$, the "arg min" above is not well-defined. If the projection $\prod$ is well-defined, we can simply apply the projection algorithm

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

# Projections

▶ Consider the projection method applied specifically to the gradient algorithm. Recall that the vector $-\nabla f(\boldsymbol{x})$ points in the direction of maximum rate of decrease of $f$ at $\boldsymbol{x}$. This was the basis for gradient methods for unconstrained optimization, which have the form $\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)})$, where $\alpha_k$ is the step size. The choice of the step size depends on the particular gradient algorithm. For example, recall that in the steepest descent algorithm, $\alpha_k = \arg\min_{\alpha \geq 0} f(\boldsymbol{x}^{(k)} - \alpha \nabla f(\boldsymbol{x}^{(k)}))$.

# Projections

▸ The projected version of the gradient algorithm has the form

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

▸ We refer to the above as the *projected gradient algorithm.*

# Example

▸ Consider the problem

$$minimize \ \frac{1}{2} x^T Q x$$
$$subject \ to \ \|x\|^2 = 1$$

where $Q = Q^T > 0$. Suppose that we apply a *fixed-step-size projected gradient algorithm* to this problem.

▸ a. Derive a formula for the update equation for the algorithm. You may assume that the argument in the projection operator to obtain $x^{(k)}$ is never zero.

# Example

▸ The projection operator in this case simply maps any vector to the closest point on the unit circle. Therefore, the projection operator is given by $\prod[x] = x/\|x\|$, provided that $x \neq 0.$ The update equation is

$$x^{(k+1)} = \beta_k\left(x^{(k)} - \alpha Q x^{(k)}\right) = \beta_k(I - \alpha Q)x^{(k)}$$

where $\beta_k = 1/\left\|(I - \alpha Q)x^{(k)}\right\|$ (i.e., it is whatever constant scaling is needed to make $x^{(k+1)}$ have unit norm).

# Example

▸ b. It is possible for the algorithm not to converge to an optimal solution even if the step size $\alpha > 0$ is taken to be arbitrarily small?

▸ If we start with $\boldsymbol{x}^{(0)}$ being an eigenvector of $\boldsymbol{Q}$, then $\boldsymbol{x}^{(k)} = \boldsymbol{x}^{(0)}$ for all $k$. Therefore, if the corresponding eigenvalue is not the smallest, then clearly the algorithm is stuck at a point that is not optimal.

# Example

▸ c. Show that for $0 < \alpha < 1/\lambda_{max}$ (where $\lambda_{max}$ is the largest eigenvalue of $\boldsymbol{Q}$), the fixed-step-size projected gradient algorithm (with step size $\alpha$) converges to an optimal solution, provided that $\boldsymbol{x}^{(0)}$ is not orthogonal to the eigenvectors of $\boldsymbol{Q}$ corresponding to the smallest eigenvalue. (Assume that the eigenvalues are distinct.)

# Example

▸ We have

$$x^{(k+1)} = \beta_k(I - \alpha Q)x^{(k)}$$

$$= \beta_k(I - \alpha Q)\left(y_1^{(k)}v_1 + \cdots + y_n^{(k)}v_n\right)$$

$$= \beta_k(y_1^{(k)}(I - \alpha Q)v_1 + \cdots + y_n^{(k)}(I - \alpha Q)v_n)$$

▸ But $(I - \alpha Q)v_i = (1 - \alpha\lambda_i)v_i$, where $\lambda_i$ is the eigenvalue corresponding to $v_i$. Hence,

$$x^{k+1} = \beta_k(y_1^{(k)}(1 - \alpha\lambda_1)v_1 + \cdots + y_n^{(k)}(1 - \alpha\lambda_n)v_n)$$

which means that $y_i^{(k+1)} = \beta_k y_i^{(k)}(1 - \alpha\lambda_i)$.

# Example

- In other words, $y_i^{(k)} = \beta^{(k)} y_i^{(0)} (1 - \alpha\lambda_i)^k$, where $\beta^{(k)} = \prod_{i=0}^{k-1} \beta_k$. We rewrite $\boldsymbol{x}^{(k)}$ as

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

- Assuming that $y_1^{(0)} \neq 0$, we obtain

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

# Example

▸ Using the fact that $\frac{1-\alpha\lambda_i}{1-\alpha\lambda_1} < 1$ (because the $\lambda_i > \lambda_1$ for

$i > 1$ and $\alpha < 1/\lambda_{max}$), we deduce that

$$\frac{y_i^{(k)}}{y_1^{(k)}} \to 0,$$

which implies that $x^{(k)} \to v_1$, as required.

# Projected Gradient Methods with Linear Constraints

▸ Consider optimization problems of the form
$$minimize \ f(\boldsymbol{x})$$
$$subject \ to \ \ \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b}$$

where $f: R^n \to R, \boldsymbol{A} \in R^{m \times n}, m < n$, rank $\boldsymbol{A} = m, \boldsymbol{b} \in R^m$

▸ We assume throughout that $f \in C^1$.

▸ The specific structure of the constraint set allows us to compute the projection operator $\prod$ using the *orthogonal projector*. Specifically, $\prod[\boldsymbol{x}]$ can be defined using the orthogonal projector matrix $\boldsymbol{P}$ given by
$$\boldsymbol{P} = \boldsymbol{I}_n - \boldsymbol{A}^T(\boldsymbol{A}\boldsymbol{A}^T)^{-1}\boldsymbol{A}$$

▸ Two important properties of the orthogonal projector $P$ that we use in this section are

  ▸ $P = P^T$

  ▸ $P^2 = P$

▸ Lemma 23.1. Let $v \in R^n$. Then, $Pv = 0$ if and only if $v \in \mathcal{R}(A^T)$. In other words, $\mathcal{N}(P) = \mathcal{R}(A^T)$. Moreover, $Av = 0$ if and only if $v \in \mathcal{R}(P)$; that is, $\mathcal{N}(A) = \mathcal{R}(P)$.

# Projected Gradient Methods with Linear Constraints

▸ Recall that in unconstrained optimization, the first-order necessary condition for a point $x^*$ to be a local minimizer is $\nabla f(x^*) = 0$.

▸ In optimization problems with equality constraints, the Lagrange condition plays the role of the first-order necessary condition.

▸ When the constraint set takes the form $\{x : Ax = b\}$, the Lagrange condition can be written as $P\nabla f(x^*) = 0$.

▸ Proposition 23.1. Let $x^* \in R^n$ be a feasible point. Then $P\nabla f(x^*) = 0$ if and only if $x^*$ satisfies the Lagrange condition.

▸ Proof. By Lemma 23.1, $P\nabla f(x^*) = 0$ if and only if we have $\nabla f(x^*) \in \mathcal{R}(A^T)$. This is equivalent to the condition that there exists $\lambda^* \in R^m$ such that $\nabla f(x^*) + A^T\lambda^* = 0$, which together with the feasibility equation $Ax = b$, constitutes the Lagrange condition.

# Projected Gradient Methods with Linear Constraints

▸ Recall that the projected gradient algorithm has the form

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

▸ For the case where the constraints are linear, it turns out that we can express the projection $\prod$ in terms of the matrix $\boldsymbol{P}$ as follows:

$$\prod[\boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)})] = \boldsymbol{x}^{(k)} - \alpha_k \boldsymbol{P} \nabla f(\boldsymbol{x}^{(k)})$$

assuming that $\boldsymbol{x}^k \in \Omega$.

# Projected Gradient Methods with Linear Constraints

▸ In our constrained optimization problem, the vector $-\nabla f(x)$ is not necessarily a feasible direction. In other words, if $x^{(k)}$ is a feasible point and we apply the algorithm $x^{(k+1)} = x^{(k)} - \alpha_k \nabla f(x^{(k)})$, then $x^{(k+1)}$ need not be feasible. This problem can be overcome by replacing $-\nabla f(x^{(k)})$ by a vector that points in a feasible direction.

▸ Note that the set of feasible directions is simply the nullspace $\mathcal{N}(A)$ of the matrix $A$. Therefore, we should first project the vector $-\nabla f(x)$ onto $\mathcal{N}(A)$. This projection is equivalent to multiplication by the matrix $P$.

# Projected Gradient Methods with Linear Constraints

▶ In summary, in the projection gradient algorithm, we update $\boldsymbol{x}^{(k)}$ according to the equation

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

# Projected Gradient Methods with Linear Constraints

▸ Proposition 23.2. In a projected gradient algorithm, if $x^{(0)}$ is feasible, then each $x^{(k)}$ is feasible; that is, for each $k \geq 0$, $Ax^{(k)} = b$.

▸ Proof. We proceed by induction. The result holds for $k = 0$ by assumption. Suppose now that $Ax^{(k)} = b$. We now show that $Ax^{(k+1)} = b$. Observe that $P\nabla f\left(x^{(k)}\right) \in \mathcal{N}(A)$. Therefore,

$$Ax^{(k+1)} = A\left(x^{(k)} - \alpha_k P\nabla f\left(x^{(k)}\right)\right)$$
$$= Ax^{(k)} - \alpha_k AP\nabla f\left(x^{(k)}\right)$$
$$= b$$

# Projected Gradient Methods with Linear Constraints

- The projected gradient algorithm updates $x^{(k)}$ in the direction of $-P\nabla f(x^{(k)})$. This vector points in the direction of maximum rate of decrease of $f$ at $x^{(k)}$ along the surface defined by $Ax = b$, as described in the following argument.

- Let $x$ be any feasible point and $d$ a feasible direction such that $\|d\| = 1$. The rate of increase of $f$ at $x$ in the direction $d$ is $\langle \nabla f(x), d \rangle$. Next, we note that because $d$ is a feasible direction, it lies in $\mathcal{N}(A)$ and hence by Lemma 23.1, we have $d \in \mathcal{R}(P) = \mathcal{R}(P^T)$. So, there exists $v$ such that $d = Pv$.

▸ Hence,
$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{d} \rangle = \langle \nabla f(\boldsymbol{x}), \boldsymbol{P}^T \boldsymbol{v} \rangle = \langle \boldsymbol{P} \nabla f(\boldsymbol{x}), \boldsymbol{v} \rangle$$

▸ By the Cauchy-Schwarz inequality,
$$\langle \boldsymbol{P} \nabla f(\boldsymbol{x}), \boldsymbol{v} \rangle \leq \|\boldsymbol{P} \nabla f(\boldsymbol{x})\| \|\boldsymbol{v}\|$$

with equality if and only if he direction of $\boldsymbol{v}$ is parallel with the direction of $\boldsymbol{P} \nabla f(\boldsymbol{x})$. Therefore, the vector $-\boldsymbol{P} \nabla f(\boldsymbol{x})$ points in the direction of maximum rate of decrease of $f$ at $\boldsymbol{x}$ among all feasible directions.

# Projected Gradient Methods with Linear Constraints

▸ Suppose that we have a starting point $x^{(0)}$, which we assume is feasible; that is, $Ax^{(0)} = b$. Consider the point $x = x^{(0)} - \alpha P\nabla f(x^{(0)})$, where $\alpha \in R$. As usual, the scalar $\alpha$ is called the step size.

▸ By the discussion above, $x$ is also a feasible point. Using a Taylor series expansion of $f$ about $x^{(0)}$ and the fact that $P = P^2 = P^T P$, we get

$$f\left(x^{(0)} - \alpha P\nabla f(x^{(0)})\right)$$
$$= f(x^{(0)}) - \alpha \nabla f(x^{(0)})^T P\nabla f(x^{(0)}) + o(\alpha)$$
$$= f(x^{(0)}) - \alpha \left\| P\nabla f(x^{(0)}) \right\|^2 + o(\alpha)$$

▶ Thus, if $\boldsymbol{P}\nabla f\left(\boldsymbol{x}^{(0)}\right) \neq 0$, that is, $\boldsymbol{x}^{(0)}$ does not satisfy the Lagrange condition, then we can choose an $\alpha$ sufficiently small such that $f(\boldsymbol{x}) < f(\boldsymbol{x}^{(0)})$, which means that $\boldsymbol{x} = \boldsymbol{x}^{(0)} - \alpha \boldsymbol{P}\nabla f(\boldsymbol{x}^{(0)})$ is an improvement over $\boldsymbol{x}^{(0)}$. This is the basis for the projected gradient algorithm $\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \boldsymbol{P}\nabla f(\boldsymbol{x}^{(k)})$, where the initial point $\boldsymbol{x}^{(0)}$ satisfies $\boldsymbol{A}\boldsymbol{x}^{(0)} = \boldsymbol{b}$ and $\alpha_k$ is some step size.

# Projected Gradient Methods with Linear Constraints

▸ A well-known variant: *projected steepest descent algorithm*, where $\alpha_k$ is given by

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\boldsymbol{x}^{(k)} - \alpha \boldsymbol{P} \nabla f(\boldsymbol{x}^{(k)}))$$

▸ Theorem 23.1. If $\{\boldsymbol{x}^{(k)}\}$ is the sequence of points generated by the projected steepest descent algorithm and if $\boldsymbol{P} \nabla f(\boldsymbol{x}^{(0)}) \neq 0$, then $f(\boldsymbol{x}^{(k+1)}) < f(\boldsymbol{x}^{(k)})$.

▸ If for some $k$, we have $P\nabla f(x^{(k)}) = 0$, then by Proposition 23.1 the point $x^{(k)}$ satisfies the Lagrange condition. This condition can be used as a stopping criterion for the algorithm. Note that if this case, $x^{(k+1)} = x^{(k)}$.

▸ For the case where $f$ is a convex function, the condition $P\nabla f(x^{(k)}) = 0$ is, in fact, equivalent to $x^{(k)}$ being a global minimizer of $f$ over the constraint set $\{x: Ax = b\}$.

# Projected Gradient Methods with Linear Constraints

▸ Proposition 23.3. The point $x^* \in R^n$ is a global minimizer of a convex function $f$ over $\{x: Ax = b\}$ if and only if $P\nabla f(x^*) = 0$.

▸ Proof. We first write $h(x) = Ax - b$. Then, the constraints can be written as $h(x) = 0$, and the problem is of the form considered in earlier chapters.

▸ Note that $Dh(x) = A$. Hence, $x^* \in R^n$ is a global minimizer of $f$ if and only if the Lagrange condition holds (see Theorem 22.8). By Proposition 23.1, this is true if and only if $P\nabla f(x^*) = 0$, and this completes the proof.

# Lagrangian Algorithm for Equality Constraints

▸ The basic idea is to use gradient algorithms to update simultaneously the decision variable and Lagrange multiplier vector.

▸ Consider the following problem
$$minimize \ f(\boldsymbol{x})$$
$$subject \ to \ \ \boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{0}$$

▸ where $\boldsymbol{h}: R^n \rightarrow R^m$. Recall that for this problem the Lagrangian function is given by
$$l(\boldsymbol{x}, \boldsymbol{\lambda}) = f(\boldsymbol{x}) + \boldsymbol{\lambda}^T \boldsymbol{h}(\boldsymbol{x})$$

▸ Assume that $f, \boldsymbol{h} \in C^2$; as usual, denote the Hessian of the Lagrangian by $\boldsymbol{L}(\boldsymbol{x}, \boldsymbol{\lambda})$

# Lagrangian Algorithm for Equality Constraints

▸ The Lagrangian algorithm for this problem is given by

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

$$\boldsymbol{\lambda}^{(k+1)} = \boldsymbol{\lambda}^{(k)} + \beta_k \boldsymbol{h}\left(\boldsymbol{x}^{(k)}\right)$$

▸ Notice that the update equation for $\boldsymbol{x}^{(k)}$ is a gradient algorithm for minimizing the Lagrangian with respect to its $\boldsymbol{x}$ argument, and the update equation for $\boldsymbol{\lambda}^{(k)}$ is a gradient algorithm for maximizing the Lagrangian with respect to its $\boldsymbol{\lambda}$ argument.

▸ Because only the gradient is used, the method is also called the *first-order Lagrangian algorithm.*

# Lagrangian Algorithm for Equality Constraints

▶ Lemma 23.2. For the Lagrangian algorithm for updating $x^{(k)}$ and $\lambda^{(k)}$, the pair $(x^*, \lambda^*)$ is a fixed point if and only if it satisfies the Lagrange condition.

▶ Below, we use $(x^*, \lambda^*)$ to denote a pair satisfying the Lagrange condition. Assume that $L(x^*, \lambda^*) > 0$. Also assume that $x^*$ is a *regular* point. With these assumptions, we are now ready to state and prove that the algorithm is locally convergent. For simplicity, we will take $\alpha_k$ and $\beta_k$ to be fixed constants (not depending on $k$), denoted $\alpha$ and $\beta$, respectively.

# Lagrangian Algorithm for Equality Constraints

▸ Theorem 23.2. For the Lagrangian algorithm for updating $x^{(k)}$ and $\lambda^{(k)}$, provided that $\alpha$ and $\beta$ are sufficiently small, there is a neighborhood of $(x^*, \lambda^*)$ such that if the pair $(x^{(0)}, \lambda^{(0)})$ is in this neighborhood, then the algorithm converges to $(x^*, \lambda^*)$ with at least a linear order of convergence.

# Lagrangian Algorithm for Inequality Constraints

▸ Consider the following optimization problem with inequality constraints:

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

where $\boldsymbol{g} \colon R^n \to R^p$.

▸ Recall that for this problem the Lagraingian function is given by

$$l(\boldsymbol{x}, \boldsymbol{\mu}) = f(\boldsymbol{x}) + \boldsymbol{\mu}^T \boldsymbol{g}(\boldsymbol{x})$$

▸ As before, assume that $f, \boldsymbol{g} \in C^2$; as usual, denote the Hessian of the Lagrangian by $\boldsymbol{L}(\boldsymbol{x}, \boldsymbol{\mu})$

# Lagrangian Algorithm for Inequality Constraints

▸ The Lagrangian algorithm for this problem is given by

$$x^{(k+1)} = x^{(k)} - \alpha_k \left( \nabla f\left(x^{(k)}\right) + Dg\left(x^{(k)}\right)^T \mu^{(k)} \right)$$

$$\mu^{(k+1)} = \left[ \mu^{(k)} + \beta_k g\left(x^{(k)}\right) \right]_+$$

where $[\cdot]_+ = \max\{\cdot, 0\}$ (applied componentwise).

▸ Notice that, as before, the update equation for $x^{(k)}$ is a gradient algorithm for minimizing the Lagrangian with respect to its $x$ argument. The update equation for $\mu^{(k)}$ is a *projected gradient algorithm* for maximizing the Lagrangian with respect to its $\mu$ argument. The reason for the projection is that the KKT multiplier vector is required to be nonnegative to satisfy the KKT condition.

# Lagrangian Algorithm for Inequality Constraints

▸ Lemma 23.3. For the Lagrangian algorithm for updating $x^{(k)}$ and $\mu^{(k)}$, the pair $(x^*, \mu^*)$ is a fixed point if and only if it satisfies the KKT condition.

▸ As before, we use the notation $(x^*, \mu^*)$ to denote a pair satisfying the KKT condition. Assume that $L(x^*, \mu^*) > 0$. Also assume that $x^*$ is a regular point. With these assumptions, we are now ready to state and prove that the algorithm is locally convergent.

# Lagrangian Algorithm for Inequality Constraints

▸ As before, we will take $\alpha_k$ and $\beta_k$ to be fixed constants (not depending on $k$), denoted $\alpha$ and $\beta$, respectively. Out analysis examines the behavior of the algorithm in two phases. In the first phase, the "nonactive" multipliers decrease to zero in finite time and remain at zero thereafter. In the second phase, the $\boldsymbol{x}^{(k)}$ iterates and the "active" multipliers converge jointly to their respective solutions, with at least a linear order of convergence.

# Lagrangian Algorithm for Inequality Constraints

▸ Theorem 23.3. For the Lagrangian algorithm for updating $x^{(k)}$ and $\mu^{(k)}$, provided that $\alpha$ and $\beta$ are sufficiently small, there is a neighborhood of $(x^*, \mu^*)$ such that if the pair $(x^{(0)}, \mu^{(0)})$ is in this neighborhood, then (1) the nonactive multipliers reduce to zero in finite time and remain at zero thereafter and (2) the algorithm converges to $(x^*, \mu^*)$ with at least a linear order of convergence.

# Penalty Methods