

A SQP-Method for Linearly Constrained Maximum Likelihood Problems

Christian Kredler

Institut für Angewandte Mathematik und Statistik
Technische Universität München

Dedicated to Prof. Klaus Ritter on his 60th birthday

Abstract. Newton-like methods are standard algorithms for unrestricted parameter estimation in a wide class of nonlinear regression models. The search directions of the here presented algorithm are solutions from a sequence of quadratic (sub-)problems (SQP) with linear constraints. The practically important generalized linear models with natural link functions, e.g. log-linear or logistic regression, lead to strictly convex optimization problems for which this easy to implement extension of Newton's method converges globally with quadratic rate. The numerical results are demonstrated at some ship damage data.

Key Words: Generalized linear models (GLM), natural link function, restricted maximum likelihood (ML-) estimation, SQP-method, global and quadratic convergence, uniform convexity.

1 Restricted ML-estimation in Generalized Linear Models

In a wide and important area of statistical data analysis we are faced with strictly convex optimization problems for which Newton-like optimization methods are more appropriate than the widely used standard algorithms, e.g. BFGS or the gradient based SQP-method of [Schittkowski (1981)]. The here proposed approach is a robust optimization tool, simply substitutable into statistical programs without using external libraries, provided a solver for quadratic programs is available; e.g. [Best & Ritter (1988)]. The convergence theorems illuminating the global and local behaviour of the suggested algorithm are proved using a technique similar to that presented in [Ritter (1982)].

The convergence properties of the investigated algorithm depend on strong convexity assumptions. Especially, these conditions are satisfied in statistical

regression problems where the canonical parameters $\boldsymbol{\theta} \in \mathbb{R}^q$ of some exponential families are involved. Consider a random variable $\mathbf{Y} \in \mathbb{R}^q$ having a density

$$h(\mathbf{y}, \boldsymbol{\theta}) = c(\mathbf{y}) \exp(\mathbf{y}^T \boldsymbol{\theta} - d(\boldsymbol{\theta})) \quad (1)$$

with respect to the Lebesgue or counting measure. The scalar valued c does not depend on $\boldsymbol{\theta}$. Note that normal, multinomial, Poisson and Gamma random variables belong to our exponential class. For the treatment of nuisance parameters like σ^2 in the normal case see e.g. [McCullagh & Nelder (1983)]. The partial derivatives d_j, d_{jk}, \dots of the analytic function d generate the cumulants of \mathbf{Y} with

$$\begin{aligned} \mathbf{E}(\mathbf{Y}) &= (d_j) \\ \mathbf{cov}(\mathbf{Y}) = \mathbf{cov}(\mathbf{Y}_j, \mathbf{Y}_k) &= (d_{jk}) \quad \text{etc.} \end{aligned} \quad (2)$$

The third order cumulants of Y can be exploited to obtain explicit bounds for variable selection techniques based on certain quadratic approximations, cf. [Kredler (1986)].

The positive definiteness of (d_{jk}) ensures the convexity properties needed for the objective functions to be analyzed below.

To clarify presentation we now restrict ourselves to the univariate case $q = 1$, although all statements apply analogously to multivariate models from which the multinomial seems to be the most important, e.g. for contingency tables and logistic discriminant analysis (see [Fahrmeir & Kredler (1984)]). Generalized linear models (GLM), i.e. regression models for exponential type random variables, have first been proposed by [Nelder & Wedderburn (1972)]. Given independent observations and covariates

$$y_t \in \mathbb{R} \quad \text{and} \quad \mathbf{z}_t \in \mathbb{R}^p, \quad t = 1, \dots, T, \quad (3)$$

where the canonical or natural parameters θ_t vary in a linear subspace

$$\theta_t = \mathbf{z}_t^T \boldsymbol{\beta} \quad (4)$$

with an unknown parameter vector $\boldsymbol{\beta} \in \mathbb{R}^p$ to be estimated. Neglecting superfluous constants the log-likelihood function for all T observations is

$$l(\boldsymbol{\beta}) = \sum_{t=1}^T \{y_t \cdot \mathbf{z}_t^T \boldsymbol{\beta} - d(\mathbf{z}_t^T \boldsymbol{\beta})\}. \quad (5)$$

According to (2)

$$-\nabla^2 l(\boldsymbol{\beta}) = \sum_{t=1}^T d''(\mathbf{z}_t^T \boldsymbol{\beta}) \mathbf{z}_t \mathbf{z}_t^T \quad \text{is positive definite} \quad (6)$$

if the matrix $(\mathbf{z}_t)_{1 \leq t \leq T}$ has full rank. Hence $-l$ is strictly convex for GLM's in the natural parameters θ_t . The existence of the ML-estimate $\hat{\boldsymbol{\beta}}$ with $l(\hat{\boldsymbol{\beta}}) \geq l(\boldsymbol{\beta})$ is not trivial, but can be ensured a priori by certain properties of the data y_t, \mathbf{z}_t , cf. [Wedderburn (1976)] and [Kaufmann (1988)]. [Albert & Anderson (1984)] guarantee especially for the binomial and multinomial case the existence of the ML-estimate $\hat{\boldsymbol{\beta}}$ under easy to verify linear separation properties of the data. For instance, $\hat{\boldsymbol{\beta}}$ cannot exist if there is a hyper-plane $H : \mathbf{0} = \mathbf{z}^T \tilde{\boldsymbol{\beta}}$ such that

$$\begin{aligned} \mathbf{z}_t^T \tilde{\boldsymbol{\beta}} &< 0 \quad \text{for all } y_t = 0 \\ \mathbf{z}_t^T \tilde{\boldsymbol{\beta}} &> 0 \quad \text{for all } y_t = 1. \end{aligned} \quad (7)$$

In section 3 a log-linear GLM for ship damage data with linear constraints for the parameter vector $\boldsymbol{\beta}$ is analyzed. This leads to

$$\max_{\boldsymbol{\beta}} \{ l(\boldsymbol{\beta}) \mid \boldsymbol{\beta} \in \mathbf{R} \subseteq \mathbb{R}^p \} \quad (8)$$

where

$$\mathbf{R} = \{ \boldsymbol{\beta} \mid \mathbf{A}\boldsymbol{\beta} \leq \mathbf{b} \}, \quad \mathbf{A} \in \mathbb{R}^{m,p}, \mathbf{b} \in \mathbb{R}^m. \quad (9)$$

In the next section we discuss the convergence properties of an SQP-algorithm to compute a linearly restricted ML-estimate of a GLM with natural link function and existing $\hat{\boldsymbol{\beta}}$.

2 SQP-Algorithm, Convergence Properties

To go along with the usual notation in the optimization literature let in this section

$$\mathbf{x} := \boldsymbol{\beta} \in \mathbb{R}^p.$$

The problem to solve is then

$$\min_{\mathbf{x}} \{ F(\mathbf{x}) \mid \mathbf{A}\mathbf{x} \leq \mathbf{b} \} \quad (10)$$

where $F(\mathbf{x}) \equiv -l(\boldsymbol{\beta})$ is twice continuously differentiable and has additionally the convex properties stated below. A linear equality constraint $\mathbf{a}^T \mathbf{x} = b$ can be replaced with two linear inequalities

$$\mathbf{a}^T \mathbf{x} \leq b \quad \text{and} \quad -\mathbf{a}^T \mathbf{x} \leq -b. \quad (11)$$

Hence, for the sake of a brief notation, the restrictions in (10) cover this case, too. However, the numerical treatment of linear equalities is much easier than that of linear inequalities, and any good numerical algorithm handles

linear equality constraints in a separate way; see also formulation of corollary 2.12. Throughout we denote gradient and Hessian of F by

$$\mathbf{g}(\mathbf{x}) := \nabla F(\mathbf{x}) \quad \text{and} \quad \mathbf{G}(\mathbf{x}) := \nabla^2 F(\mathbf{x}).$$

Def. 2.1 (Uniform Convexity)

$F \in C^2(\mathbb{R}^p)$ is said to be **uniformly convex** on a convex set $\mathbf{D} \subseteq \mathbb{R}^p$, if and only if there are $0 < \mu \leq \eta$ such that

$$\mu \|\mathbf{x}\|^2 \leq \mathbf{x}^T \mathbf{G}(\mathbf{u}) \mathbf{x} \leq \eta \|\mathbf{x}\|^2 \quad (12)$$

for all $\mathbf{x} \in \mathbb{R}^p$ and all $\mathbf{u} \in \mathbf{D}$.

Uniformly convex functions are strictly convex and attain their unique minimum. $F(x) = x^4$ shows that the reverse is not true.

However, for certain applications like natural link function GLM's the following theorem gives a concrete characterization depending on the existence of a minimizer $\bar{\mathbf{x}}$. As mentioned in the previous section this can be done, for instance in practically relevant GLM's, by an a priori inspection of the data involved in the nonlinear function F .

Theorem 2.2

If $F \in C^2(\mathbb{R}^p)$ has a positive definite Hessian for all $\mathbf{x} \in \mathbb{R}^p$ and attains its minimum at some $\bar{\mathbf{x}}$, then F is **uniformly convex** on the level sets $\mathbf{N}_\alpha := \{\mathbf{x} \in \mathbb{R}^p \mid F(\mathbf{x}) \leq F(\bar{\mathbf{x}})\}$ for all $\alpha \in \mathbb{R}$.

Proof:

We first show that the level sets are bounded. The compactness follows, because they are closed. Pick $\omega > 0$, let $\alpha := F(\bar{\mathbf{x}}) + \omega$, and consider $\mathbf{N}_\alpha := \{\mathbf{x} \mid F(\mathbf{x}) \leq \alpha\}$, which is convex since F is strictly convex. If it was not bounded there would exist a $\mathbf{s} \neq \mathbf{0} \in \mathbb{R}^p$ and $\varphi(\sigma) := F(\bar{\mathbf{x}} - \sigma \mathbf{s}) \leq \alpha$ for all $\sigma \geq 0$. φ is strictly convex, and hence $F(\bar{\mathbf{x}}) = \varphi(0) < \varphi(\sigma)$, for all $\sigma \neq 0$. With $\delta := \varphi(1) - \varphi(0) > 0$ we obtain for all $\sigma > 1$

$$\varphi(0) + \delta = \varphi(1) = \varphi\left(\frac{1}{\sigma}\sigma + \left(1 - \frac{1}{\sigma}\right)0\right) < \frac{1}{\sigma}\varphi(\sigma) + \left(1 - \frac{1}{\sigma}\right)\varphi(0),$$

which implies $\varphi(0) + \delta < \frac{1}{\sigma}(\varphi(\sigma) - \varphi(0)) + \varphi(0)$, and further

$$0 < \delta < \frac{1}{\sigma}(\varphi(\sigma) - \varphi(0)) \leq \frac{1}{\sigma}(\alpha - \varphi(0)) = \frac{1}{\sigma}(\omega + \varphi(0) - \varphi(0)) = \frac{\omega}{\sigma}$$

for all $\sigma > 1$, which contradicts $\delta = \varphi(1) - \varphi(0) > 0$. Now, the extremal eigenvalues attain their minimum and maximum on the compact set \mathbf{N}_α . \square

The theorem covers most important generalized linear models including binomial, multinomial and Poisson variables. For the Γ -distribution and related cases the theorem can be extended, with some technical caution, to $F \in C^2(\mathbf{D})$ where \mathbf{D} is an appropriate, convex subset of \mathbb{R}^p . In the proofs of the following statements we need only the uniform convexity of F on $\mathbf{N}_\circ \cap \mathbf{R}$, the intersection of the level set and the feasible region $\mathbf{R} = \{\mathbf{x} \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$. Trivially, this property can always be obtained by restricting F with positive definite Hessian to a box $\mathbf{v} \leq \mathbf{x} \leq \mathbf{w}$.

2.1 A General Descent Method

We consider an iterative algorithm which, starting from $\mathbf{x}_0 \in \mathbb{R}^p$, computes for $j = 0, 1, 2, \dots$ a search direction \mathbf{s}_j and a stepsize σ_j defining the new point

$$\mathbf{x}_{j+1} := \mathbf{x}_j - \sigma_j \mathbf{s}_j,$$

where $\mathbf{g}_j^T \mathbf{s}_j > 0$ for $\mathbf{g}_j := \mathbf{g}(\mathbf{x}_j)$.

Def. 2.3 (Armijo-Goldstein Rule)

Choose $0 < \delta < \frac{1}{2}$. In \mathbf{x}_j find the smallest integer $\nu_j \geq 0$ such that

$$F(\mathbf{x}_j) - F(\mathbf{x}_j - (\frac{1}{2})^{\nu_j} \mathbf{s}_j) \geq \delta \cdot (\frac{1}{2})^{\nu_j} \cdot \mathbf{g}_j^T \mathbf{s}_j. \quad (13)$$

Finally, let $\sigma_j = (\frac{1}{2})^{\nu_j}$.

Def. 2.4 (Quadratic Approximation)

With $\mathbf{g}_j := \mathbf{g}(\mathbf{x}_j)$, $\mathbf{G}_j := \mathbf{G}(\mathbf{x}_j)$ the quadratic Taylor approximation is $F(\mathbf{x}_j - \mathbf{s}) \approx F(\mathbf{x}_j) + Q_j(\mathbf{s})$ where

$$Q_j(\mathbf{s}) := -\mathbf{g}_j^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{G}_j \mathbf{s}. \quad (14)$$

Let throughout $\boldsymbol{\lambda} := (\lambda_1, \dots, \lambda_m)^T \leq \mathbf{0}$ and

$$\mathbf{A} := \begin{pmatrix} \mathbf{a}_1^T \\ \cdots \\ \mathbf{a}_m^T \end{pmatrix}, \quad \mathbf{b} := \begin{pmatrix} b_1 \\ \cdots \\ b_m \end{pmatrix}.$$

Def. 2.5 (Kuhn-Tucker Conditions, Lagrange Multipliers)

The tuple $(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}) \in \mathbb{R}^{p+m}$, $\bar{\boldsymbol{\lambda}} \leq \mathbf{0}$, is called **Kuhn-Tucker pair** of problem (10) if the following conditions hold

$$\mathbf{g}(\bar{\mathbf{x}}) - \mathbf{A}^T \bar{\boldsymbol{\lambda}} = \mathbf{0}, \quad \bar{\boldsymbol{\lambda}} \leq \mathbf{0}, \quad (15)$$

$$\mathbf{A} \bar{\mathbf{x}} \leq \mathbf{b}, \quad (16)$$

$$\bar{\lambda}_k (\mathbf{a}_k^T \bar{\mathbf{x}} - b_k) = 0, \quad 1 \leq k \leq m. \quad (17)$$

$\bar{\mathbf{x}}$ satisfying the above conditions is called **stationary point** and $\bar{\boldsymbol{\lambda}}$ is the vector of optimal **Lagrange multipliers**.

In our case F is convex and differentiable, hence each stationary point $\bar{\mathbf{x}}$ is a global optimum of (10).

Algorithm SQPStep 0: Initialization

Choose a feasible \mathbf{x}_0 , i.e. $\mathbf{A} \mathbf{x}_0 \leq \mathbf{b}$, $0 < \delta < \frac{1}{2}$.
Let $j := 0$, and goto Step 1.

Step 1: Search Direction, Stopping Criterion

Compute $\mathbf{g}_j := \mathbf{g}(\mathbf{x}_j)$, $\mathbf{G}_j := \mathbf{G}(\mathbf{x}_j)$ and

$$\mathbf{s}_j := \arg \min_{\mathbf{s}} \{ Q_j(\mathbf{s}) \mid \mathbf{A}(\mathbf{x}_j - \mathbf{s}) \leq \mathbf{b} \} \quad (18)$$

If $\mathbf{s}_j = \mathbf{0}$, then **STOP**, else goto Step 2.

Step 2: Stepsize

Choose σ_j according to Armijo-Goldstein.
Goto Step 3.

Step 3: Update:

Define

$$\mathbf{x}_{j+1} := \mathbf{x}_j - \sigma_j \mathbf{s}_j,$$

replace j with $j + 1$, and goto Step 1.

Theorem 2.6 *If F is uniformly convex on $N_o \cap \mathbf{R}$, then algorithm SQP is well defined. The generated iterates $\{\mathbf{x}_j\}$ ensure $F(\mathbf{x}_{j+1}) < F(\mathbf{x}_j)$. If $\mathbf{s}_j = \mathbf{0}$, then \mathbf{x}_j is a stationary point of (10).*

Proof:

The subproblem (18) is strictly convex, hence \mathbf{s}_j is unique. We now have to show that $-\mathbf{s}_j$ is a descent direction, and that $\mathbf{s}_j = \mathbf{0}$ can only occur in a stationary point.

$\nabla Q_j(\mathbf{s}) = -\mathbf{g}_j + \mathbf{G}_j \mathbf{s}$. The Kuhn-Tucker conditions for (18) imply

$$-\mathbf{g}_j + \mathbf{G}_j \mathbf{s}_j + \mathbf{A}^T \boldsymbol{\lambda}_j = \mathbf{0}; \quad \boldsymbol{\lambda}_j \leq \mathbf{0}, \quad (19)$$

$$\mathbf{A}(\mathbf{x}_j - \mathbf{s}_j) \leq \mathbf{b}, \quad (20)$$

$$\boldsymbol{\lambda}_j^T (\mathbf{A}(\mathbf{x}_j - \mathbf{s}_j) - \mathbf{b}) = \mathbf{0}. \quad (21)$$

(21) follows from the usual complementarity condition. (19) yields

$$\mathbf{s}_j = \mathbf{G}_j^{-1}(\mathbf{g}_j - \mathbf{A}^T \boldsymbol{\lambda}_j) \quad (22)$$

which is well defined since \mathbf{G}_j is positive definite by assumption. With $\boldsymbol{\lambda}_j \leq \mathbf{0}$ from (19) and $\mathbf{A} \mathbf{x}_j - \mathbf{b} \leq \mathbf{0}$ we obtain by the complementarity condition (21)

$$\boldsymbol{\lambda}_j^T \mathbf{A} \mathbf{s}_j \geq 0. \quad (23)$$

Further, (19) multiplied by \mathbf{s}_j gives

$$\mathbf{g}_j^T \mathbf{s}_j = \mathbf{s}_j^T \mathbf{G}_j \mathbf{s}_j + \boldsymbol{\lambda}_j^T \mathbf{A} \mathbf{s}_j. \quad (24)$$

(23) and the uniform convexity yield finally

$$\mathbf{g}_j^T \mathbf{s}_j \geq \mu \|\mathbf{s}_j\|^2 > 0. \quad (25)$$

This guarantees the descent property of the algorithm. We remark further that \mathbf{x}_{j+1} stays within the feasible region $\mathbf{R} = \{\mathbf{x} \mid \mathbf{A} \mathbf{x} \leq \mathbf{b}\}$ because $\mathbf{x}_j \in \mathbf{R}$, $\mathbf{x}_j - \mathbf{s}_j \in \mathbf{R}$ and $\sigma_j \in (0, 1]$. In case $\mathbf{s}_j = \mathbf{0}$ we conclude from (19) that \mathbf{x}_j must be a stationary point of (10), and hence is optimal due to the convexity of F . \square

Def. 2.7 (Lagrangian)

$$L(\mathbf{x}, \boldsymbol{\lambda}) := F(\mathbf{x}) - \boldsymbol{\lambda}^T (\mathbf{A} \mathbf{x} - \mathbf{b}), \quad \boldsymbol{\lambda} \leq \mathbf{0}, \quad (26)$$

denotes the **Lagrangian function** of problem (10).

Theorem 2.8 (Global Convergence)

Provided F is uniformly convex on $N_o \cap \mathbf{R}$, the intersection of level set and feasible region, the SQP-algorithm converges to a stationary point $\bar{\mathbf{x}}$ of (10) with

$$\lim_{j \rightarrow \infty} \nabla L(\mathbf{x}_j) = \mathbf{0}. \quad (27)$$

Proof:

We follow the scheme of [Ritter (1982)]. To get a contradiction assume that $\|\nabla L(\mathbf{x}_j)\| = \|\mathbf{g}_j - \mathbf{A}^T \boldsymbol{\lambda}_j\| > \varepsilon > 0$ for all $j \in J$ where J is an infinite subset of \mathbf{N} . Hence by (22)

$$\|\mathbf{s}_j\| \geq \frac{1}{\eta} \|\mathbf{g}_j - \mathbf{A}^T \boldsymbol{\lambda}_j\| \geq \frac{\varepsilon}{\eta}, \quad \text{for all } j \in J. \quad (28)$$

Further, by (25)

$$\frac{\mathbf{g}_j^T \mathbf{s}_j}{\|\mathbf{s}_j\|} \geq \mu \|\mathbf{s}_j\| \geq \frac{\varepsilon \mu}{\eta}, \quad \text{for all } j \in J. \quad (29)$$

According to Armijo-Goldstein we choose σ_j such that

$$a_j(\sigma_j) := \frac{F(\mathbf{x}_j) - F(\mathbf{x}_j - \sigma_j \mathbf{s}_j)}{\sigma_j \mathbf{g}_j^T \mathbf{s}_j} \geq \delta. \quad (30)$$

Taylor expansion with $\mathbf{u}_j \in [\mathbf{x}_j, \mathbf{x}_j - \sigma_j \mathbf{s}_j]$ gives

$$\begin{aligned} a_j(\sigma) &\geq 1 - \frac{\|\mathbf{g}(\mathbf{u}_j) - \mathbf{g}_j\| \cdot \|\mathbf{s}_j\|}{\mathbf{g}_j^T \mathbf{s}_j} \\ &= 1 - \frac{\|\mathbf{g}(\mathbf{u}_j) - \mathbf{g}_j\|}{\frac{\mathbf{g}_j^T \mathbf{s}_j}{\|\mathbf{s}_j\|}} \geq 1 - \frac{\eta}{\varepsilon \mu} \|\mathbf{g}(\mathbf{u}_j) - \mathbf{g}_j\|. \end{aligned} \quad (31)$$

As $\mathbf{g}(\cdot)$ is uniformly continuous on the compact set N_o , there is a $\tau > 0$ such that $\|\mathbf{g}(\mathbf{u}_j) - \mathbf{g}_j\| \leq (\varepsilon \mu)/(2\eta)$ for $\|\sigma \mathbf{s}_j\| \leq \tau$. Then (31) with $\|\sigma \mathbf{s}_j\| \leq \tau$ yields

$$a_j(\sigma) \geq 1 - \frac{1}{2} = \frac{1}{2} > \delta, \quad \text{for all } j \in J,$$

which means that (13) is satisfied. In the Armijo-Goldstein rule we choose the smallest ν_j (i.e. $\sigma_j = \frac{1}{2}^{\nu_j}$ as large as possible). Hence for $j \in J$ with $\sigma_j \leq 1$

$$\|\sigma_j \mathbf{s}_j\| = \left\| \left(\frac{1}{2}\right)^{\nu_j} \mathbf{s}_j \right\| \geq \min\{\|\mathbf{s}_j\|, \frac{\tau}{2}\},$$

and further for all $j \in J$

$$\begin{aligned} F(\mathbf{x}_j) - F(\mathbf{x}_j - \sigma_j \mathbf{s}_j) &\geq \delta \sigma_j \mathbf{g}_j^T \mathbf{s}_j = \delta \|\sigma_j \mathbf{s}_j\| \frac{\mathbf{g}_j^T \mathbf{s}_j}{\|\mathbf{s}_j\|} \geq \delta \min\{\|\mathbf{s}_j\|, \frac{\tau}{2}\} \frac{\varepsilon \mu}{\eta} \\ &\geq \delta \min\{\frac{\varepsilon}{\eta}, \frac{\tau}{2}\} \frac{\varepsilon \mu}{\eta} > 0. \end{aligned} \quad (32)$$

Note that the bound is independent of j and \mathbf{x}_j . As J is infinite F cannot be bounded below, which contradicts the compactness of \mathbf{N}_o . The remaining Kuhn-Tucker conditions are obtained from those holding for subproblem (18).

F is uniformly, hence strictly convex. Therefore, each stationary point $\bar{\mathbf{x}}$ is a global minimum of the linearly constrained problem (10). Further, $\bar{\mathbf{x}}$ is unique such that the sequence $\{\mathbf{x}_j\}$ converges to $\bar{\mathbf{x}}$. \square

2.2 Convergence Rate

A specific property of quadratically convergent Newton-like algorithms is that the stepsize σ_j converges to 1. This is also the case for our SQP-algorithm.

Lemma 2.9 (Stepsize $\sigma_j = 1$)

If F is uniformly convex on $\mathbf{N}_o \cap \mathbf{R}$ and $\{\mathbf{s}_j\}, \{\sigma_j\}$ are generated by algorithm SQP there is a $j_0 \in \mathbf{N}$ such that $\sigma_j = 1$ for $j \geq j_0$.

Proof:

According to Taylor there is some $\mathbf{u}_j \in [\mathbf{x}_j, \mathbf{x}_j - \mathbf{s}_j]$ such that

$$F(\mathbf{x}_j - \mathbf{s}_j) = F(\mathbf{x}_j) - \mathbf{g}_j^T \mathbf{s}_j + \frac{1}{2} \mathbf{s}_j^T \mathbf{G}(\mathbf{u}_j) \mathbf{s}_j. \quad (33)$$

We have to show $a_j(1) \geq \delta > 0$ for the Armijo-Goldstein function from (30).

$$\begin{aligned} a_j(1) &= \frac{F(\mathbf{x}_j) - F(\mathbf{x}_j - 1 \cdot \mathbf{s}_j)}{1 \cdot \mathbf{g}_j^T \mathbf{s}_j} \\ &= \frac{1}{\mathbf{g}_j^T \mathbf{s}_j} \left(\frac{1}{2} \mathbf{g}_j^T \mathbf{s}_j + \frac{1}{2} \mathbf{g}_j^T \mathbf{s}_j - \frac{1}{2} \mathbf{s}_j^T (\mathbf{G}_j + \mathbf{G}(\mathbf{u}_j) - \mathbf{G}_j) \mathbf{s}_j \right). \end{aligned} \quad (34)$$

With $\boldsymbol{\lambda}_j^T \mathbf{A} \mathbf{s}_j \geq 0$ from (23) we obtain

$$a_j(1) \geq \frac{1}{2} + \frac{\mathbf{g}_j^T \mathbf{s}_j - \mathbf{s}_j^T \mathbf{G}_j \mathbf{s}_j - \boldsymbol{\lambda}_j^T \mathbf{A} \mathbf{s}_j}{2 \mathbf{g}_j^T \mathbf{s}_j} - \frac{\mathbf{s}_j^T (\mathbf{G}(\mathbf{u}_j) - \mathbf{G}_j) \mathbf{s}_j}{2 \mathbf{g}_j^T \mathbf{s}_j}. \quad (35)$$

According to (19) the second term on the right hand side vanishes. Hence with (25)

$$\begin{aligned}
a_j(1) &\geq \frac{1}{2} - \frac{1}{2} \frac{\mathbf{s}_j^T (\mathbf{G}(\mathbf{u}_j) - \mathbf{G}_j) \mathbf{s}_j}{\mathbf{g}_j^T \mathbf{s}_j} \\
&\geq \frac{1}{2} - \frac{1}{2} \frac{\|\mathbf{G}(\mathbf{u}_j) - \mathbf{G}_j\| \cdot \|\mathbf{s}_j\|^2}{\mu \|\mathbf{s}_j\|^2} \\
&= \frac{1}{2} - \frac{1}{2} \frac{\|\mathbf{G}(\mathbf{u}_j) - \mathbf{G}_j\|}{\mu}
\end{aligned} \tag{36}$$

which converges to $\frac{1}{2} > \delta$ since $\lim \mathbf{G}(\mathbf{u}_j) = \lim \mathbf{G}_j = \mathbf{G}(\bar{\mathbf{x}})$. \square

So far the only used regularity condition was the uniform convexity of F . For technical reasons we now need assumptions on the linear independence of the gradients \mathbf{a}_k corresponding to the active constraints. Denote by $\bar{K} := \{1 \leq k \leq m \mid \mathbf{a}_k^T \bar{\mathbf{x}} = b_k\}$ the index set of active constraints at $\bar{\mathbf{x}}$, and by $\bar{\boldsymbol{\lambda}} := (\bar{\lambda}_1, \dots, \bar{\lambda}_m)^T$ the optimal Lagrange multipliers.

Def. 2.10 (Assumption A)

The optimization problem (10) is said to satisfy **Assumption A** if and only if F is twice continuously differentiable and if $(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}})$, $\bar{\boldsymbol{\lambda}} \leq \mathbf{0}$, is a Kuhn-Tucker pair of (10) according to Def. 2.5 and

$$\mathbf{s}^T \mathbf{G}(\bar{\mathbf{x}}) \mathbf{s} > 0 \text{ for all } \mathbf{s} \neq \mathbf{0} \text{ with } \mathbf{a}_k^T \mathbf{s} = 0, \quad k \in \bar{K}, \tag{37}$$

and

$$\bar{\lambda}_k < 0 \text{ for } k \in \bar{K}, \quad (\text{strict complementary slackness}) \tag{38}$$

and if

$$\bar{\mathbf{A}} := (\mathbf{a}_k)_{k \in \bar{K}}^T \in \mathbb{R}^{r,p} \text{ has full rank.} \tag{39}$$

Clearly, (37) follows for the special models of section 1 directly from the uniform convexity of F on the level set N_c . The gradients of the active constraints in (39) can only be linearly independent if linear equalities are stated separately, and not in the way proposed for theoretical purposes in (11). Without loss of generality the corresponding Lagrange multiplier can be chosen positive; otherwise we redefine the equation by $-\mathbf{a}_k^T \mathbf{x} = -b_k$.

With stepsize $\sigma_j = 1$ our algorithm SQP is nothing else but the method of [Wilson (1963)] for linear constraints; see also [Gill et al. (1989)]. Among

others the quadratic convergence properties of Wilson's method are investigated in [Robinson (1974)]. Here, we only mention those aspects of the proof which can be exploited algorithmically. The subproblems (18) are viewed as perturbations of the original problem for which, under certain conditions, the solution converges to that of (10).

Theorem 2.11 (Active Constraints, Lagrange Multipliers)

Let the $\{\mathbf{x}_j\}$ generated by algorithm SQP converge to $\bar{\mathbf{x}}$ and λ_j be the Lagrange multipliers of the quadratic subproblems (18). If (10) satisfies **Assumption A** at $(\bar{\mathbf{x}}, \bar{\lambda})$ then the Lagrange multipliers λ_j of the quadratic subproblems (18) converge to the optimal multipliers $\bar{\lambda}$. Further, there is an j_0 such that

$$(\lambda_j)_k < 0 \quad \text{if} \quad \bar{\lambda}_k < 0 \quad \text{and} \quad (40)$$

$$\mathbf{a}_k^T \mathbf{x}_j < b_k \quad \text{if} \quad \mathbf{a}_k^T \bar{\mathbf{x}} < b_k \quad (41)$$

for all $j \geq j_0$.

Proof: See theorem 2.1 of [Robinson (1974)], which is an application of the implicit function theorem. \square

Hence we can be sure to pick ultimately the correct active constraints. In such a case, however, the solution of the subproblem (18) can be found very efficiently.

Corollary 2.12 (Newton Direction in Subspaces)

Choose $\mathbf{S} \in \mathbb{R}^{p-r,p}$ with orthonormal columns and $\bar{\mathbf{A}}^T \mathbf{S} = \mathbf{0}$, where according to (39), $r = |\bar{K}| = \text{rank}(\bar{\mathbf{A}})$. Under the assumptions of theorem 2.11 there is an index j_0 such that the subproblem (18) can be replaced with

$$\mathbf{w}_j := \arg \min_{\mathbf{w}} \{ F(\mathbf{x}_j - \mathbf{S} \mathbf{w}) \mid \mathbf{w} \in \mathbb{R}^{p-r} \}. \quad (42)$$

$\mathbf{s}_j := \mathbf{S} \mathbf{w}_j$ of (18) can explicitly be computed from

$$\mathbf{S}^T \mathbf{G}_j \mathbf{S} \mathbf{w}_j = \mathbf{S}^T \mathbf{g}_j. \quad (43)$$

Theorem 2.13 (Quadratic Convergence)

Under the assumptions of theorem 2.11 the algorithm SQP converges quadratically to the optimal point $\bar{\mathbf{x}}$ if $\nabla^2 F$ satisfies a Lipschitz condition in a neighbourhood of $\bar{\mathbf{x}}$.

Proof:

Consider \mathbf{x}_{j_0} with $\bar{\mathbf{A}}\mathbf{x}_{j_0} = (b_k)_{k \in \bar{K}}$ and

$$\bar{F}(\mathbf{w}) := F(\mathbf{x}_{j_0} - \mathbf{S}\mathbf{w}), \quad (44)$$

being uniformly convex with

$$\nabla^2 \bar{F}(\mathbf{w}) = \mathbf{S}^T \mathbf{G}(\mathbf{x}_{j_0} - \mathbf{S}\mathbf{w}) \mathbf{S}. \quad (45)$$

With $\nabla^2 F$ also $\nabla^2 \bar{F}$ satisfies a Lipschitz condition. Hence Newton's method converges quadratically. According to theorem 2.11 and corollary 2.12 in a neighbourhood close enough to $\bar{\mathbf{x}}$ the sequences generated by Newton's method for \bar{F} and by algorithm SQP coincide. This completes the proof. \square

Corollary 2.14 (Unconstrained Case)

All the stated global and local convergence properties remain valid for the case without linear restrictions. This fact is well known and has been exploited since ever to compute unconstrained ML-estimates in generalized linear models (GLM's) efficiently.

Remark:

1. The convergence results can be extended to problems with nonlinear constraints, cf. e.g. [Robinson (1974)]. From the data analysis point of view distribution results are hard to obtain even for linear constraints (see [McDonald & Diamond (1990)], [Piegorisch (1990)] and [Nyquist (1991)]). This may explain, why this paper is restricted to the linearly constrained case, where convergence results and algorithms can be formulated in quite a simple manner. Nevertheless, a review of the proofs shows that the assumptions can be weakened. The identities (19) - (25) play a crucial role throughout, and there may be scenarios where (23) and (25) hold on $N_0 \cap \mathbf{R}$ and F is far from being convex.
2. General link functions for GLM's are of considerable practical interest. Like for unconstrained ML-estimates the Hessian should then be replaced with the positive (semi-)definite Fisher scoring matrix $\tilde{\mathbf{G}}_j$, cf. e.g. [Fahrmeir & Kredler (1984)], which is of the same approximation type as the Gauß-Newton-Matrix in Nonlinear Regression. The resulting modified SQP-algorithm can be expected to have satisfactory numerical behaviour. Weaker convergence properties than those stated here can be obtained under appropriate conditions, e.g. work with $\tilde{\mathbf{G}}_j + \kappa \mathbf{I}$ instead of $\tilde{\mathbf{G}}_j$, if necessary, to obtain good enough descent directions.

3 A GLM with Parameter Restrictions

Ship Damage Data (Lloyd's Register of Shipping)				
cf. [McCullagh & Nelder (1983)], p. 137				
Ship type	Year of construction	Period of operation	aggregate months service	number of damage incidents
A	1960-64	1960-74	127	0
A	1960-64	1975-79	63	0
A	1965-69	1960-74	1095	3
A	1965-69	1975-79	1095	4
A	1970-74	1960-74	1512	6
A	1970-74	1975-79	3353	18
A	1975-79	1960-74	0	0*
A	1975-79	1975-79	2244	11
B	1960-64	1960-74	44882	39
B	1960-64	1975-79	17176	29
B	1965-69	1960-74	28609	58
B	1965-69	1975-79	20370	53
B	1970-74	1960-74	7064	12
B	1970-74	1975-79	13099	44
B	1975-79	1960-74	0	0*
B	1975-79	1975-79	7117	18
C	1960-64	1960-74	1179	1
C	1960-64	1975-79	552	1
C	1965-69	1960-74	781	0
C	1965-69	1975-79	676	1
C	1970-74	1960-74	783	6
C	1970-74	1975-79	1948	2
C	1975-79	1960-74	0	0*
C	1975-79	1975-79	274	1
D	1960-64	1960-74	251	0
D	1960-64	1975-79	105	0
D	1965-69	1960-74	288	0
D	1965-69	1975-79	192	0
D	1970-74	1960-74	349	2
D	1970-74	1975-79	1208	11
D	1975-79	1960-74	0	0*
D	1975-79	1975-79	2051	4
E	1960-64	1960-74	45	0
E	1960-64	1975-79	0	0*
E	1965-69	1960-74	789	7
E	1965-69	1975-79	437	7
E	1970-74	1960-74	1157	5
E	1970-74	1975-79	2161	12
E	1975-79	1960-74	0	0*
E	1975-79	1975-79	542	1

* Necessarily empty cells for ships not (yet) in action.

The company has to balance financial risk and is interested in a statistical analysis that allows a distinction between categories of safe ships and others.

Damage Rates in ‰				
	year of construction			
ship type	60-64	65-69	70-74	75-79
A	0.0	3.2	4.9	4.9
B	1.1	2.3	2.3	2.5
C	1.2	0.7	2.9	3.6
D	0.0	0.0	8.3	2.0
E	0.0	11.4	5.1	1.8

Lowest risk is observed for ship types B and C and highest for type E. Unexpectedly, the oldest ships appear to be the safest, whereas those built in 1965-1974 seem to yield highest risk. To validate this conjecture we define an appropriate generalized linear model (GLM, see section 1) satisfying all assumptions made in the previous section. The ML-estimation of the model parameters was done with the SQP-algorithm described in the previous section. The reported numerical results can be found in [Klinger (1992)].

$T = 40$ design vectors

$$\mathbf{z} = (1, z_{.1}, \dots, z_{.8}, z_{.9})^T$$

are defined with effect coding for the factors. For instance, reference category 60 – 64 yields

$(z_{.1}, z_{.2}, z_{.3})$	year of construction
$(1, 0, 0)$	60 – 64
$(0, 1, 0)$	65 – 69
$(0, 0, 1)$	70 – 74
$(-1, -1, -1)$	75 – 79

Reference category ship type E defines $z_{.4} \cdots z_{.7}$ for the remaining 4 categories A,B,C,D. $z_{.8} = 1$ is for period of operation in 75-79, whereas $z_{.8} = -1$ means 60-74. Finally, $z_{.9}$ describes the monthly service hours of the individual ship.

We follow [McCullagh & Nelder (1983)] and view the number y_t of damage incidents as independent Poisson variables. This gives a log-linear GLM with parameter vector $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_9)^T$

$$\ln(E[y_t]) = \beta_0 + \sum_{j=1}^9 z_{.j} \beta_j, \quad (46)$$

which yields the log-likelihood function

$$l(\boldsymbol{\beta}) = \sum_{i=1}^N [y_i \mathbf{z}_i^T \boldsymbol{\beta} - \exp(\mathbf{z}_i^T \boldsymbol{\beta})]. \quad (47)$$

The analysis of interactions is omitted here. Unrestricted ML-estimation gives

$$\hat{\boldsymbol{\beta}} = (-5.175, \quad -0.448, \quad 0.214, \quad 0.312, \quad \dots, \quad 0.905)$$

β_0 $\underbrace{60-64}$ $\underbrace{65-69}$ $\underbrace{70-74}$ \dots service hours
 year of construction (75-79 is reference category)

Due to effect coding the reference category parameter becomes

$$\hat{\beta}_{75-79} = -\sum_{j=1}^3 \hat{\beta}_j = -0.078. \quad (48)$$

Again, the conjecture that older ships are safer than new ones can be expressed in terms of the parameters of (46). Statistically spoken, we want to know whether the parameters corresponding to earlier construction periods are significantly smaller than the others. This can be formulated in the following hypotheses

$$H_0 : \beta_1 = \beta_2 = \beta_3$$

versus

$$H_1 : \beta_1 \leq \beta_2 \leq \beta_3.$$

$\hat{\boldsymbol{\beta}}$ satisfies restriction H_1 . Under the equality constraints of H_0 we get

$$\tilde{\boldsymbol{\beta}} = (-3.940, \quad 0.084, \quad 0.084, \quad 0.084, \quad \dots, \quad 0.759)$$

β_0 $\underbrace{65-69}$ $\underbrace{70-74}$ $\underbrace{75-79}$ \dots service hours
 year of construction

Following the theory of [Wollak (1987)], [McDonald & Diamond (1990)], and [Fahrmeir & Klinger (1994)] the corrected likelihood ratio statistic lq^* is distributed like a mixture of χ^2 -variables with specific weights. [Klinger (1992)] reports $lq^* = 19.28$ with a "p-value" of 0.0001. Hence H_1 is accepted, and the year of construction has an highly significant influence on the number of damage incidents.

One possibility to check the significance of the negative influence of β_{75-79} is to recode the variables for the construction period with new reference category 60 – 64. We obtain

$$\hat{\beta} = (-5.175, \quad 0.214, \quad 0.312, \quad -0.078, \quad \dots, \quad 0.905)$$

$$\begin{array}{ccccccc} \beta_0 & \underbrace{65-69} & \underbrace{70-74} & \underbrace{75-79} & \dots & \text{service hours} & \\ & \text{year of construction} & & & & & \text{(60-64 is reference category)} \end{array}$$

The statistical hypothesis to test is now

$$H_0: \beta_1 \leq \beta_2 \leq \beta_3$$

versus

$$H_1: \beta_1, \beta_2, \beta_3 \text{ arbitrary (not in } H_0).$$

With

$$\mathbf{A} = \begin{pmatrix} 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

H_0 is equivalent with $\mathbf{A}\beta \leq \mathbf{b}$. Under these restrictions we obtain with the algorithm of section 2 the ML-estimate

$$\tilde{\beta} = (-5.366, \quad 0.159, \quad 0.180, \quad 0.180, \quad \dots, \quad 0.935)$$

$$\begin{array}{ccccccc} \beta_0 & \underbrace{65-69} & \underbrace{70-74} & \underbrace{75-79} & \dots & \text{service hours} & \\ & \text{year of construction} & & & & & \end{array}$$

Again, according to [Klinger (1992)], the corrected likelihood ratio statistic is $lq^* = 3.04$ with a p-value of 0.0804. So hypothesis H_0 is accepted, which finally confirms the conjecture that in the considered data set older ships are more robust than newer ones.

Examples with linear parameter restrictions for binomial credit scoring and endodontic risk data are discussed in [Fahrmeir & Klinger (1994)].

4 Conclusion

SQP-methods are appropriate for a wider class of parameter restricted estimation problems than the here described generalized linear models with canonical parameters. If analytic first derivatives are available the FORTRAN-algorithm of [Schittkowski (1981)] will usually produce good results.

Sometimes the exact model is not known a priori and alternatives are to be checked before the best model is found. In this case the derivative coding for each new model is not only time consuming but also a source of many errors. Numerical difference quotients instead of analytic derivatives are an

obvious but frequently not numerically stable alternative. In that situation automatic differentiation is a good way out. The PC-program PADMOS, cf. [Greiner, Kredler & Wagenpfeil (1992)], has been developed especially for such purposes. Text files contain the Pascal-like description of log-likelihood function and constraints. In the data file for each observation an extra line is reserved. A comfortable user interface with built-in editor enables a convenient input and modification of ML-problems. Models with some hundred observations and up to 15 parameters can be analyzed. Box, linear and even nonlinear type constraints are accepted. The execution time of PADMOS is about 10 times of that needed by analytic derivative procedures. But for any chosen model first **and** second derivatives are computed automatically, and hence the broad variety of powerful optimization algorithms, including methods with directions of negative curvature, is applicable, and ensures good convergence properties.

Acknowledgement

I wish to thank Joachim Klinger and Stefan Wagenpfeil for helpful comments and discussions.

References

- [Albert & Anderson (1984)] Albert A. and Anderson J.A.: On the existence of maximum likelihood estimates in logistic regression models. *Biometrika* **71**, 1-10.
- [Best & Ritter (1988)] Best M.J. and Ritter K.: A quadratic programming algorithm. *ZOR* **32**, 271-297.
- [Fahrmeir & Klinger (1994)] Fahrmeir L. and Klinger J.: Estimating and testing generalized linear models under inequality restrictions. *Statistical Papers* **35**, 211-229.
- [Fahrmeir & Kredler (1984)] Fahrmeir L. und Kredler Ch.: Verallgemeinerte lineare Modelle. In: Fahrmeir L. und Hamerle A. (Hrsg.): *Multivariate statistische Verfahren*. De Gruyter, Berlin.
- [Gill et al. (1989)] Gill P.E., Murray W., Saunders M.A. and Wright M.H.: Constrained nonlinear programming. In: G.L. Nemhauser et al. (eds.): *Handbooks in Operations Research and Management Science*. North-Holland, Amsterdam.
- [Greiner, Kredler & Wagenpfeil (1992)] Greiner M., Kredler Ch. and Wagenpfeil St.: *Nonlinear optimization with PADMOS: User's guide and algorithms*. Report **366**, DFG Schwerpunkt: Anwendungsbezogene Optimierung und Steuerung, TU München.

- [Kaufmann (1988)] H.: On existence and uniqueness of maximum likelihood estimates in quantal and ordinal response models. *Metrika* **35**, 291-313.
- [Klinger (1992)] J.: Tests für Ungleichungsrestriktionen in generalisierten linearen Modellen. Diplomarbeit (Referent. Prof. L. Fahrmeir), LMU München.
- [Kredler (1986)] Ch.: Behaviour of third order terms in quadratic approximations of LR-statistics in multivariate generalized linear models. *The Annals of Statistics* **14**, 326-335.
- [McCullagh & Nelder (1983)] McCullagh P. and Nelder J.A.: Generalized linear models. Chapman and Hall, London.
- [McDonald & Diamond (1990)] McDonald P. and Diamond I.: On the fitting of generalized linear models with nonnegativity parameter constraints. *Biometrics* **46**, 201-206.
- [Nelder & Wedderburn (1972)] Nelder J. and Wedderburn R.W.M.: Generalized linear models. *J. Roy. Statist. Soc. A* **135**, 370-384.
- [Nyquist (1991)] H.: Restricted estimation of generalized linear models. *Applied Statistics*, **40**, 133-141.
- [Piegorisch (1990)] W.: One-side-significance tests for generalized linear models under dichotomous response. *Biometrics* **46**, 309-316.
- [Ritter (1982)] K.: Numerical methods for nonlinear programming problems. In: Korte B. (ed.): MODERN APPLIED MATHEMATICS — Optimization and Operations Research. North-Holland Publ. Comp. Amsterdam.
- [Robinson (1974)] S.M.: Perturbed Kuhn-Tucker points and rates of convergence for a class of nonlinear programming algorithms. *Mathematical Programming* **7**,1-16.
- [Schittkowski (1981)] K.: The nonlinear programming method of Wilson, Han and Powell with an augmented Lagrangian type line search function. Part 1: Convergence analysis. *Numer. Math.* **38**, 83-114.
- [Wedderburn (1976)] R.W.M.: On the existence and uniqueness of the maximum likelihood estimates for certain generalized linear models. *Biometrika* **63**, 27-32.
- [Wilson (1963)] R.B.: A simplicial algorithm for concave programming. PhD thesis, Harvard University, Cambridge.
- [Wollak (1987)] F.A.: An exakt test for multiple inequality and equality constraints in the linear model. *JASA* **82**, 782-793.