

# Time-lag in Derivative Convergence for Fixed Point Iterations\*

Andreas Griewank and Andrej Ponomarenko

## Abstract

In an earlier study it was proven and experimentally confirmed on a 2D Euler code that fixed point iterations can be differentiated to yield first and second order derivatives of implicit functions that are defined by state equations. It was also asserted that the resulting approximations for reduced gradients and Hessians converge with the same R-factor as the underlying fixed point iteration.

A closer look reveals now that nevertheless these derivative values lag behind the functions values in that the ratios of the corresponding errors grow proportional to the iteration counter or its square towards infinity. This rather subtle effect is caused mathematically by the occurrence of nontrivial Jordan blocks associated with degenerate eigenvalues. We elaborate the theory and report its confirmation through numerical experiments.

## Résumé

Une étude antérieure a prouvé et vérifié expérimentalement sur un code Euler 2D que les calculs itératifs avec point fixe peuvent être différenciés pour obtenir les dérivées aux premier et deuxième ordres des fonctions implicites définies par des équations d'état. On considérait également que les itérées correspondantes des gradients et Hessiens réduits convergent à la même vitesse que l'itération de point fixe d'origine.

Cette étude plus détaillée révèle néanmoins que ces dérivées convergent avec un certain retard par rapport aux valeurs de la fonction. En effet le rapport des erreurs correspondantes croît vers l'infini proportionnellement au compteur d'itérations ou à son carré. Mathématiquement, cet effet plutôt subtil est causé par l'apparition de blocs de Jordan correspondant à des valeurs propres dégénérées. Nous construisons un modèle théorique de cet effet et nous le validons par des expériences numériques.

## 1 Introduction and Assumption

The effect to be analyzed arises in the context of design optimization by what has been called piggy-back optimization [2]. Design optimization problems are distinguished from general non-linear programming problems (NLP) by the fact that the vector of variables  $x$  is a priori partitioned into a state vector  $y \in Y$  and a set of design variables  $u \in U$ . For application of this scenario in computational fluidynamics see for example [7], [5], [6], and [4]. Throughout we assume that the "user" has provided an iteration function

$$G : Y \times U \rightarrow Y$$

that is contractive with respect to an inner product norm on  $Y$  so that for all  $u \in U$  and  $y, \tilde{y} \in Y$

$$\|G(y, u) - G(\tilde{y}, u)\| \leq \rho \|y - \tilde{y}\|.$$

\*Supported by the DFG research center "Mathematics for the key technologies" (FZT 80) in Berlin.



Here  $\varrho < 1$  may vary continuously as a function of the design  $u$  and its exact size will usually not be available to a practical algorithm.

As an immediate consequence it follows by the Banach fixed point theorem that for fixed  $u$  and any initial  $y_0 \in Y$  the sequence  $\{y_k\}$  generated by

$$y_{k+1} = G(y_k, u)$$

must converge to the unique fixed point  $y_* = y_*(u)$  with  $y_* = G(y_*, u)$ . In other words, the assumptions made so far ensure that one can obtain for any  $u$  a solution  $y_*(u)$ , a process which one may call "simulate" the underlying system. In a practical simulation the variables  $u$  and  $y$  will often be restricted to open subsets of the spaces  $U$  and  $Y$ , respectively.

In order to progress from simulation to design we require more smoothness of  $G$ , namely, that it is at least once continuously differentiable in the joint variable vector  $(y, u)$ . The same assumption will be made for the objective function

$$f : Y \times U \rightarrow \mathbb{R},$$

which is meant to be minimized. Provided at least  $f \in C^1(Y, U)$ , one can obtain in a completely automated fashion the adjoint iteration function

$$\tilde{G}(y, \tilde{y}, u) \equiv \tilde{y} G_y(y, u) + f_y(y, u). \quad (1)$$

Here subscripts denote partial differentiation and  $\tilde{y}$  like the gradient  $f_y$  is considered a row-vector belonging to the dual space of  $Y$ , which we identify with the Hilbert space  $Y$  itself. Then we have in the induced matrix and operator norm

$$\varrho(u) = \max_{y \in Y} \|G_y(y, u)\| \leq \varrho < 1$$

so that also in the dual norm

$$\|\tilde{G}(y, \tilde{y}, u) - \tilde{G}(y, \hat{\tilde{y}}, u)\| \leq \varrho \|\tilde{y} - \hat{\tilde{y}}\|$$

for any two row-vectors  $\tilde{y}, \hat{\tilde{y}} \in \tilde{Y} \equiv Y$ .

## 2 Piggy-Back Convergence of Adjoint

Throughout the remainder of this paper we consider  $u$  as constant and may therefore omit it occasionally as an argument in analyzing the simultaneous iteration

$$\begin{bmatrix} y_{k+1} \\ \tilde{y}_{k+1} \end{bmatrix} = \begin{bmatrix} G(y_k, u) \\ \tilde{G}(y_k, \tilde{y}_k, u) \end{bmatrix} \quad (2)$$

Even if  $G$  is merely  $C^1$  and thus  $G_y(y) = G_y(y, u)$  continuous with respect to  $y$  it follows from  $y_k \rightarrow y_*$  that  $G_y(y_k, u) \rightarrow G_y(y_*, u)$  and hence the adjoint iterates  $\tilde{y}_k$  converge to  $\tilde{y}_*$  the unique solution of the adjoint equation

$$\tilde{y}_* = \tilde{y}_* G_y(y_*, u) + f_y(y_*, u) \quad (3)$$

The vector  $\tilde{y}_*$  can be used to compute the so called reduced gradient

$$\tilde{u}_* = \tilde{y}_* G_u(y_*, u) + f_u(y_*, u) \quad (4)$$





This row vector represents the total derivatives of  $f$  with respect to  $u$ , after the elimination of the state vector  $y$  using the implicit function theorem. In order to be more specific about the rate of convergence we assume that  $G_y$  and  $f_y$  are Lipschitz continuous with respect to  $y$  so that for some  $\nu > 0$

$$\|G_y(\bar{y}, u) - G_y(y, u)\| \leq \nu \|\bar{y} - y\| \geq \|f_y(\bar{y}, u) - f_y(y, u)\|.$$

Then we find immediately for the discrepancies  $\Delta y_k = y_k - y_*$  and  $\Delta \bar{y}_k = \bar{y}_k - \bar{y}_*$  that

$$\begin{aligned} \|\Delta \bar{y}_{k+1}\| &\leq \|\bar{y}_k G_y(y_k, u) - \bar{y}_* G_y(y_*, u)\| + \|f_y(y_k, u) - f_y(y_*, u)\| \\ &= \|\Delta \bar{y}_k G_y(y_k, u) + \bar{y}_* (G_y(y_k, u) - G_y(y_*, u))\| + \nu \|y_k - y_0\| \\ &\leq \varrho \|\Delta \bar{y}_k\| + \|\Delta y_k\| (\|\bar{y}_*\| + 1) \nu \end{aligned}$$

Consequently we have for any weighted error combination

$$z_k = \|\Delta y_k\| + \omega \|\Delta \bar{y}_k\|$$

the recurrence

$$\begin{aligned} z_{k+1} &\leq \varrho \|\Delta y_k\| + \omega (\varrho \|\Delta \bar{y}_k\| + \nu (\|\bar{y}_*\| + 1) \|\Delta y_k\|) \\ &= (\varrho + \omega \nu (\|\bar{y}_*\| + 1)) \|\Delta y_k\| + \omega \varrho \|\Delta \bar{y}_k\| \\ &\leq (\varrho + \omega \nu (\|\bar{y}_*\| + 1)) z_k \end{aligned}$$

This implies for any  $\omega < (1 - \varrho) / (\nu (\|\bar{y}_*\| + 1))$  the Q-linear convergence result

$$\limsup_{k \rightarrow \infty} z_{k+1} / z_k \leq \varrho + \omega \nu (\|\bar{y}_*\| + 1) < 1$$

By standard arguments one derives the R-linear convergence results

$$\limsup_{k \rightarrow \infty} \sqrt[k]{\|\Delta y_k\|} \leq \lim_{k \rightarrow \infty} \sqrt[k]{z_k / \omega} \leq \varrho + \omega \nu (\|\bar{y}_*\| + 1) < 1$$

Taking the infimum over all  $\omega > 0$  one finally obtains as in [1]

$$\limsup_{k \rightarrow \infty} \sqrt[k]{\|\Delta y_k\|} \leq \varrho \geq \limsup_{k \rightarrow \infty} \sqrt[k]{\|\Delta \bar{y}_k\|}.$$

Here the inequality on the right was just added for comparison. Since these convergence speed cannot be improved under our assumptions (namely  $G_y$  has maximal norm  $\varrho$  and is Lipschitz continuous with respect to  $y$ ) one may arrive at the conclusion that the sequences  $\{y_k\}$  and  $\{\bar{y}_k\}$  converge essentially at the same speed. In fact this claim has been made repeatedly in the literature and the present author has suffered from the same impression for a long time. On the other hand there has been the persistent notion that the convergence of derivatives is lagging behind those of the underlying fixed point iterates.

### 3 Relative Convergence Speed of First Adjoints

In the remainder of this paper we require that  $Y = \mathbb{R}^n$  and  $U = \mathbb{R}^m$  are finite dimensional Euclidean spaces so that all linear operators can be identified with their matrix presentation. Assuming furthermore, that  $G$  and  $f$  are twice Lipschitz-continuously differentiable, we may rewrite the recurrence (2) as

$$\begin{bmatrix} y_{k+1} \\ \bar{y}_{k+1} \end{bmatrix} = \begin{bmatrix} G(y_k, u) \\ N_y(y_k, \bar{y}_k, u) \end{bmatrix} \quad (5)$$





Here we have expressed the  $G$  from (3) as the gradient of the function

$$N(y, \bar{y}, u) \equiv \bar{y} G(y, u) + f(y, u)$$

with respect to  $y$ . Notice that this function  $N$  differs from the familiar Lagrange function  $L$  of the optimization problem  $\text{Min}(f(y, u))$  s.t.  $G(y, u) - \bar{y} = 0$  by the shift

$$\bar{y} y = N(y, \bar{y}, u) - L(y, \bar{y}, u)$$

Consequently, we have

$$N_y = L_y + \bar{y} \quad \text{and} \quad N_{\bar{y}} = L_{\bar{y}} + y \quad \text{but} \quad N_u = L_u$$

and, for the subsequent analysis more importantly, all second derivatives are identical:

$$N_{yy} = L_{yy}, \quad N_{y\bar{y}} = L_{y\bar{y}}, \quad N_{\bar{y}\bar{y}} = L_{\bar{y}\bar{y}}$$

Differentiating (5) we obtain the block-triangular Jacobian

$$J_k \equiv \frac{\partial(y_{k+1}, \bar{y}_{k+1})}{\partial(y_k, \bar{y}_k)} = \begin{bmatrix} G_y(y_k, u) & 0 \\ N_{yy}(y_k, \bar{y}_k, u) & G_y^T(y_k, u) \end{bmatrix}$$

Obviously we have the characteristic polynomial

$$\det(J_k - \lambda I) = \det^2(G_y(y_k, u) - \lambda I)$$

which means that  $J_k$  has the same eigenvalues as  $G_y(y_k, u)$  but all of them with the multiplicity 2. Moreover, at least one of these eigenvalues will be defective and thus generate a Jordan block of dimension greater than 2, unless the quadratic form  $v^T N_{yy} v$  vanishes for all eigenvectors  $v$  of  $G_y(y_k, u)$ . As a consequence one can deduce a linear-geometric decline in the adjoint error as follows.

Linearizing about the fixed point  $(y_*, \bar{y}_*)$  we obtain the Taylor expansion

$$\begin{bmatrix} \Delta y_{k+1} \\ \Delta \bar{y}_{k+1} \end{bmatrix} = \begin{bmatrix} A & 0 \\ B & A^T \end{bmatrix} \begin{bmatrix} \Delta y_k \\ \Delta \bar{y}_k \end{bmatrix} + O(\|\Delta y_k\|^2 + \|\Delta \bar{y}_k\|^2)$$

where  $A \equiv G_y(y_*, u)$  and  $B \equiv N_{yy}(y_*, \bar{y}_*, u)$ . From this it follows by induction using the B-linear convergence of  $\|\Delta y_k\| + \|\Delta \bar{y}_k\|$  that for any  $k$  and  $j > 0$

$$\begin{bmatrix} \Delta y_{k+j} \\ \Delta \bar{y}_{k+j} \end{bmatrix} = \begin{bmatrix} A & 0 \\ B & A^T \end{bmatrix}^j \begin{bmatrix} \Delta y_k \\ \Delta \bar{y}_k \end{bmatrix} + O(\|\Delta y_k\|^2 + \|\Delta \bar{y}_k\|^2) \tag{6}$$

Similarly it can be easily verified by induction that

$$J_k^j \equiv \begin{bmatrix} A & 0 \\ B & A^T \end{bmatrix}^j = \begin{bmatrix} A^j & 0 \\ \sum_{i=1}^j (A^T)^{i-1} B A^{j-i} & (A^T)^j \end{bmatrix} \tag{7}$$

To simplify the matrix on the bottom left we assume that  $A = G_y$  is real diagonalizable so that

$$A = T \Gamma T^{-1} \quad \text{with} \quad \Gamma = \text{diag}\{\gamma_j\}_{j=1}^n$$

where

$$\rho_* = \max_{1 \leq j \leq n} |\gamma_j| \leq \rho < 1$$





Then we can perform a two stage reduction to obtain the Jordan-like representation

$$\begin{aligned} \begin{bmatrix} A & 0 \\ B & A^T \end{bmatrix} &= \begin{bmatrix} T & 0 \\ 0 & T^{-T} \end{bmatrix} \begin{bmatrix} \Gamma & 0 \\ T^T B T & \Gamma \end{bmatrix} \begin{bmatrix} T^{-1} & 0 \\ 0 & T^T \end{bmatrix} \\ &= \begin{bmatrix} T & 0 \\ 0 & T^{-T} \end{bmatrix} \begin{bmatrix} J & 0 \\ C^T & I \end{bmatrix} \begin{bmatrix} \Gamma & 0 \\ D & \Gamma \end{bmatrix} \begin{bmatrix} I & 0 \\ C & I \end{bmatrix} \begin{bmatrix} T^{-1} & 0 \\ 0 & T^T \end{bmatrix} \end{aligned} \tag{8}$$

Here  $D$  is the (real) diagonal of  $T^T B T$  and  $C = -C^T$  is the antisymmetric solution of the Liapunov equation

$$\Gamma C - C \Gamma = T^T B T - D.$$

It is well known that the linear mapping from  $C$  to  $\Gamma C - C \Gamma$  has the  $\alpha^2$  eigenvalues  $\gamma_i - \gamma_j$  and the eigenvectors  $e_i e_j^T$  for  $1 \leq i, j \leq n$ , so that the Liapunov equation must be solvable if all eigenvalues of  $A$  are distinct, as we will assume for simplicity. Then it follows immediately that the  $j$ -th power of  $J$ , is given by

$$\begin{bmatrix} A & 0 \\ B & A^T \end{bmatrix}^j = \begin{bmatrix} T & 0 \\ (C T^{-1})^j & T^{-j} \end{bmatrix} \begin{bmatrix} \Gamma^j & 0 \\ D \Gamma^{j-1} & \Gamma^j \end{bmatrix} \begin{bmatrix} T^{-1} & 0 \\ C T^{-1} & T^T \end{bmatrix} \tag{9}$$

Thus we see that unless the diagonal  $D$  of  $T^T B T$  happens to vanish there might be a pretty strong growth in the adjoint error component  $\Delta \hat{y}_k$ . To compare it to the original error  $\Delta y_k$  itself we firstly have to analyze its recurrence a bit more carefully. Using the Lipschitz constant  $\nu$  one finds by standard estimates

$$\|\Delta y_{k+1} - A \Delta y_k\| \leq \nu \|\Delta y_k\|^2.$$

Abbreviating

$$\Delta \hat{y}_{k+1} = T^{-1} \Delta y_{k+1} \quad \text{and} \quad \kappa = \|T\| \|T^{-1}\|$$

one has consequently

$$\|\Delta \hat{y}_{k+1} - \Gamma \Delta \hat{y}_k\| \leq \nu \kappa \|\Delta \hat{y}_k\|^2.$$

Using this estimate and the assumption that there is only one eigenvalue, say  $\gamma_1$  with the maximal modulus  $\rho_* = |\gamma_1|$  one can show that the angle between  $\hat{y}_k$  and the first Cartesian basis vector  $e_1$  satisfies a recurrence that has exactly one stable fixed point namely 0. Thus we have generically

$$\lim_{k \rightarrow \infty} \frac{|e_1^T \Delta \hat{y}_k|}{\|\Delta \hat{y}_k\|} = 1 \quad \text{and} \quad \lim_{k \rightarrow \infty} \frac{\|\Delta \hat{y}_k\|}{\rho_*^k} = c_* \in (0, \infty).$$

From (6) and (9) it then follows that

$$T^T \Delta \hat{y}_{k+j} = [C^T \Gamma^j + j D \Gamma^{j-1} + \Gamma^j C] \Delta \hat{y}_k + \Gamma^j T^T \Delta y_k + O(\|\Delta y_k\|^2 + \|\Delta \hat{y}_k\|^2)$$

Assuming finally that the first diagonal element  $d_1$  of  $D$  does not vanish we find that the second term in the above expansion dominates as  $j$  grows and we get the limit exactly

$$\lim_{k \rightarrow \infty} \frac{1}{k} \frac{\|T^T \Delta \hat{y}_k\|}{\|T^{-1} \Delta y_k\|} = \frac{d_1}{\rho_*} \tag{10}$$

so that obviously

$$\lim_{k \rightarrow \infty} \frac{1}{k} \frac{\|\Delta \hat{y}_k\|}{\|\Delta y_k\|} \leq \kappa \frac{d_1}{\rho_*}$$





and approximately

$$\frac{\|\Delta \hat{y}_k\|}{\|\Delta y_k\|} \sim k$$

Hence we see that the convergence of the adjoint vectors  $\hat{y}_k$  really lags behind that of the underlying iterates  $y_k$  even though both sequences have the same R-factor  $\rho_k$ .

### 4 Convergence of Second Order Adjoints

The above analysis can be extended to second derivatives representing products of the projected Hessian with certain direction vectors. More specifically, after picking a direction  $\hat{u} \in U$  we may append (2) by the iterations

$$\hat{y}_{k+1} \equiv \hat{G}(y_k, \hat{y}_k, u, \hat{u}) \equiv G_y(y_k, u)\hat{y}_k + G_u(y_k, u)\hat{u} \tag{11}$$

and

$$\begin{aligned} \hat{\hat{y}}_{k+1} \equiv \hat{\hat{G}}(y_k, \hat{y}_k, \hat{y}_k, \hat{y}_k, u, \hat{u}) &\equiv \hat{y}_k G_y + \hat{y}_k G_{yy} \hat{y}_k + f_{yu} \hat{y}_k + \hat{y}_k G_{yu} \hat{u} + f_{yu} \hat{u} \\ &= \hat{y}_k G_y(y_k, u) + N_{yy}(y_k, \hat{y}_k, u)\hat{y}_k + N_{yu}(y_k, \hat{y}_k, u)\hat{u} \end{aligned} \tag{12}$$

where all derivatives of  $G$  and  $f$  are evaluated at the current argument  $(y_k, u)$ . Then an analysis along the lines of Section 3 shows that the  $\hat{y}_k$  and  $\hat{\hat{y}}_k$  also converge R-linearly to respective fixed points  $\hat{y}_*$  and  $\hat{\hat{y}}_*$  solving

$$\hat{y}_* = \hat{G}(y_*, \hat{y}_*, u, \hat{u}) \quad \text{and} \quad \hat{\hat{y}}_* = \hat{\hat{G}}(y_*, \hat{y}_*, \hat{y}_*, \hat{y}_*, u, \hat{u})$$

The vector  $\hat{y}_*$  represents the feasible direction in state space associated with the variation  $\hat{u}$  in the design space. The vector  $\hat{\hat{y}}_*$  can be used to compute

$$\hat{\hat{u}}_* \equiv \hat{\hat{y}}_* G_u(y_*, u) + N_{uy}(y_*, \hat{y}_*, u)\hat{y}_* + N_{uu}(y_*, \hat{y}_*, u)\hat{u} \tag{13}$$

which represents the product of the reduced Hessian with the direction  $\hat{u}$ . To analyze the speed of convergence more carefully let us consider the extended Jacobian

$$\frac{\partial(y_{k+1}, \hat{y}_{k+1}, \hat{\hat{y}}_{k+1}, \hat{y}_{k+1})}{\partial(y_k, \hat{y}_k, \hat{y}_k, \hat{y}_k)} = \begin{bmatrix} G_y(y_k, u) & 0 & 0 & 0 \\ N_{yy}(y_k, \hat{y}_k, u) & G_y^T(y_k, u) & 0 & 0 \\ P(y_k, \hat{y}_k, u, \hat{u}) & 0 & G_u(y_k, u) & 0 \\ H(y_k, \hat{y}_k, \hat{y}_k, \hat{y}_k, u, \hat{u}) & P(y_k, \hat{y}_k, u, \hat{u})^T & N_{yu}^T(y_k, \hat{y}_k, u) & G_u^T(y_k, u) \end{bmatrix}$$

where

$$\begin{aligned} P(y, \hat{y}, u, \hat{u}) &= G_{yy}(y, u)\hat{y} + G_{yu}(y, u)\hat{u} \\ H(y, \hat{y}, \hat{y}, \hat{y}, u, \hat{u}) &= \hat{y} G_{yy}(y, u) + N_{yy}(y, \hat{y}, u)\hat{y} + N_{yu}(y, \hat{y}, u)\hat{u} \end{aligned}$$

We notice that the matrix  $H$  is symmetric, while  $P$  is general and the values of these two square matrices at the fixed point  $(y_*, \hat{y}_*, \hat{y}_*, \hat{y}_*)$  are independent of each other as well as  $A$  and  $B$ .

We are looking now for estimates of the corresponding discrepancies  $\Delta \hat{y}_k \equiv y_k - \hat{y}_*$  and  $\Delta \hat{\hat{y}}_k \equiv \hat{\hat{y}}_k - \hat{\hat{y}}_*$  in addition to the  $\Delta y_k$  and  $\Delta \hat{y}_k$  considered before. Similarly to (6) we obtain the linearization

$$\begin{bmatrix} \Delta y_{k+1} \\ \Delta \hat{y}_{k+1} \\ \Delta \hat{\hat{y}}_{k+1} \\ \Delta \hat{y}_{k+1} \end{bmatrix} = \begin{bmatrix} A & 0 & 0 & 0 \\ B & A^T & 0 & 0 \\ P & 0 & A & 0 \\ H & P^T & B & A^T \end{bmatrix}^T \begin{bmatrix} \Delta y_k \\ \Delta \hat{y}_k \\ \Delta \hat{\hat{y}}_k \\ \Delta \hat{y}_k \end{bmatrix} + O(\|\Delta y_k\|^2 + \|\Delta \hat{y}_k\|^2 + \|\Delta \hat{\hat{y}}_k\|^2 + \|\Delta \hat{y}_k\|^2)$$





Using the same transformation as in (8) the  $j$ -th power can be rewritten as following

$$\begin{bmatrix} A & 0 & 0 & 0 \\ B & A^T & 0 & 0 \\ P & 0 & A & 0 \\ H & P^T & B & A^T \end{bmatrix}^j = \begin{bmatrix} T & 0 & 0 & 0 \\ T^{-T}C^T & T^{-T} & 0 & 0 \\ 0 & 0 & T & 0 \\ 0 & 0 & T^{-T}C^T & T^{-T} \end{bmatrix} \begin{bmatrix} \Gamma^j & 0 & 0 & 0 \\ jD\Gamma^{j-1} & \Gamma^j & 0 & 0 \\ \hat{P}_j & 0 & \Gamma^j & 0 \\ H_j & \hat{P}_j^T & j\Gamma^{j-1}D & \Gamma^j \end{bmatrix} \begin{bmatrix} T^{-1} & 0 & 0 & 0 \\ CT^{-1} & T^T & 0 & 0 \\ 0 & 0 & T^{-1} & 0 \\ 0 & 0 & CT^{-1} & T^T \end{bmatrix}$$

where with  $\hat{P} \equiv T^{-1}PT$  and  $\hat{H} \equiv T^THT$

$$\begin{bmatrix} \hat{P}_j & 0 \\ H_j & \hat{P}_j^T \end{bmatrix} = \sum_{i=1}^j \begin{bmatrix} \Gamma^{i-1} & 0 \\ (i-1)\Gamma^{i-2}D & \Gamma^{i-1} \end{bmatrix} \begin{bmatrix} \hat{P} & 0 \\ \hat{H} & \hat{P}^T \end{bmatrix} \begin{bmatrix} \Gamma^{j-i} & 0 \\ (j-i)\Gamma^{j-i-1}D & \Gamma^{j-i} \end{bmatrix}$$

Here we used the relation (7) once again. Hence we have the expressions

$$\begin{aligned} \hat{P}_j &= \sum_{i=1}^j \Gamma^{i-1} \hat{P} \Gamma^{j-i} \\ H_j &= \sum_{i=1}^j (i-1)\Gamma^{i-2}D\hat{P}\Gamma^{j-i} + \Gamma^{j-1}\hat{H}\Gamma^{j-i} + (j-i)\Gamma^{i-1}P\Gamma^{j-i-1}D \end{aligned}$$

Taking norms we obtain for constants  $c_1$  and  $c_2$

$$\|\hat{P}_j\| \leq c_1 j \rho_k^{j-1} \quad \text{and} \quad \|H_j\| \leq c_2 j^2 \rho_k^{j-2}$$

The later inequality is true because  $\Gamma$  has like  $A$  the spectral radius  $\rho_k < 1$ . Thus we can estimate all four error components as follows,

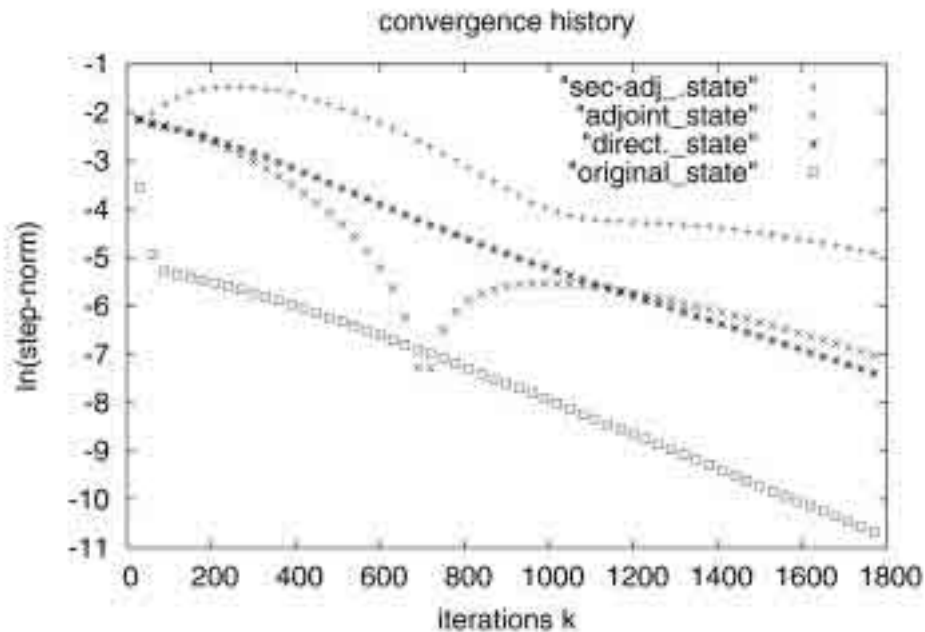
$$\begin{aligned} \|\Delta \hat{y}_{k+j}\| &\leq \rho_k^j c_{14} \|\Delta \hat{y}_k\| + O(\|\Delta \hat{y}_k\|^2) \\ \|\Delta \hat{y}_{k+j}\| &\leq \rho_k^j [c_{22} \|\Delta \hat{y}_k\| + c_{21} j \|\Delta \hat{y}_k\|] + O(\|\Delta \hat{y}_k\|^2 + \|\Delta \hat{y}_k\|^2) \\ \|\Delta \hat{y}_{k+j}\| &\leq \rho_k^j [c_{33} \|\Delta \hat{y}_k\| + c_{31} j \|\Delta \hat{y}_k\|] + O(\|\Delta \hat{y}_k\|^2 + \|\Delta \hat{y}_k\|^2) \\ \|\Delta \hat{y}_{k+j}\| &\leq \rho_k^j [c_{44} \|\Delta \hat{y}_k\| + c_{41} j^2 \|\Delta \hat{y}_k\| + c_{42} j \|\Delta \hat{y}_k\| + \|\Delta \hat{y}_k\|] \\ &\quad + O(\|\Delta \hat{y}_k\|^2 + \|\Delta \hat{y}_k\|^2 + \|\Delta \hat{y}_k\|^2 + \|\Delta \hat{y}_k\|^2) \end{aligned}$$

Using the assumption (10) one can actually obtain the proportionality relations

$$\|\Delta \hat{y}_k\| \sim k \|\Delta y_k\| \sim k \rho_k^k \quad \text{and} \quad \|\Delta \hat{\hat{y}}_k\| \sim k^2 \|\Delta y_k\| \sim k^2 \rho_k^k$$

This means in particular that the second derivatives lag behind the first derivatives by a factor of order  $k$  and thus behind the original iteration by a factor of order  $k^2$ . A closer analysis of the propagation constants in the above system of bounds might allow an optimized decision of when to start propagating first and second derivatives. If this is done too early the error  $\|\Delta y_j\|$  might still be so large that no reduction in the derivative errors occurs at all.





## 5 Numerical Results

The following results were obtained on the boundary control problem

$$\Delta_x y(x) + e^{\sigma(x)} = 0 \quad \text{for } x = (x_1, x_2) \in [0, 1]^2$$

with the periodic and Dirichlet boundary conditions

$$y(0, \zeta) = y(1, \zeta), \quad y(\zeta, 0) = \sin(2\pi\zeta), \quad y(\zeta, 1) = u(\zeta) \quad \text{for } \zeta \in [0, 1]$$

The function  $u$  is viewed as a boundary control that can be varied to minimize the objective function

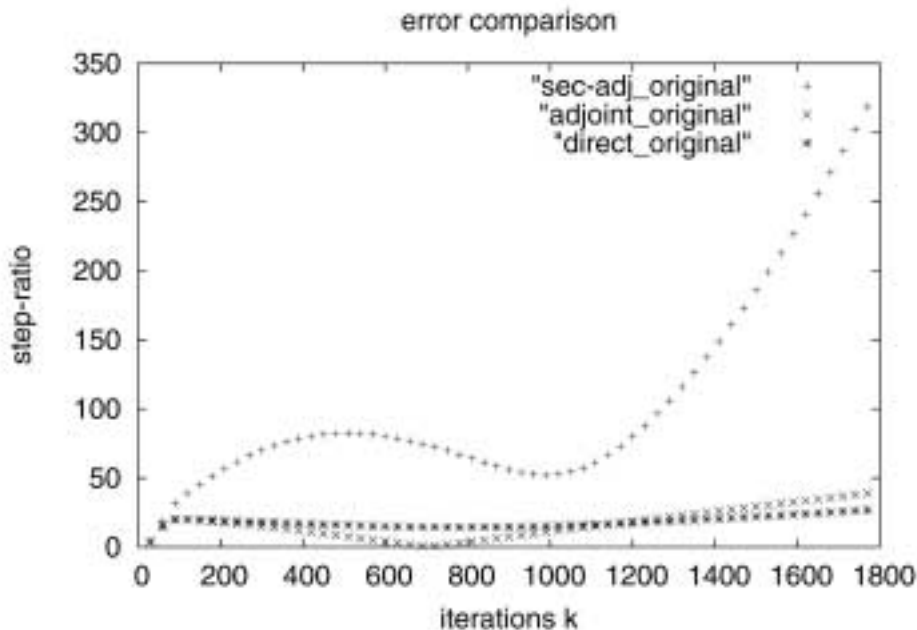
$$J(y, u) = \int_0^1 \left[ \left. \frac{\partial y(\eta, \zeta)}{\partial \eta} \right|_{\eta=0} - 4 - \cos(2\pi\zeta) \right]^2 d\zeta + \sigma \int_0^1 [u(\zeta)^2 + u'(\zeta)^2] d\zeta$$

In the following calculations we used  $\sigma = 0.001$  and set constantly  $u(\zeta) = 2.2$ . This value is not all that far from the fold point where solutions cease to exist.

We use a central difference discretization with the mesh-width  $1/12.0$  so that the resulting algebraic system involves 144 equations in as many variables. Since the nonlinearities occur only on the diagonal one can easily implement Jacobi's method to obtain the basic function  $G(y, u)$ . For this simple example we also coded by hand the corresponding derived functions  $\hat{G}$ ,  $\hat{G}$  and even  $\hat{G}$  as defined in (1, 11) and (12), respectively. The results were later confirmed using the automatic differentiation tool ADOL-C [3].

As can be seen in Fig.1 the convergence of the Jacobi method is rather slow with the common R-factor being about  $(1 - 1/300)$ . The lowest curve represents the natural logarithm of the Euclidean norm ratios  $\|y_{k+1} - y_k\| / \|y_1 - y_0\|$ , which provide some indication of the norm ratios





$\|\Delta y_k\|/\|\Delta y_0\|$ . In view of the very slow convergence this relation need certainly not be very close. Nevertheless the theory is basically confirmed with the first direct and adjoint derivatives  $\|\hat{y}_{k+1} - \hat{y}_k\|/\|\hat{y}_1 - \hat{y}_0\|$  and  $\|\hat{y}_{k+1} - \hat{y}_k\|/\|\hat{y}_1 - \hat{y}_0\|$  lagging somewhat behind and the second derivatives  $\|\hat{y}_{k+1} - \hat{y}_k\|/\|\hat{y}_1 - \hat{y}_0\|$  coming in last. The ratio between these derivative quantities and the original iterates themselves is plotted in Fig. 2. After an initial transition phase one sees quite clearly a growth proportional to  $k$  and  $k^2$  for the first and second derivatives, respectively. While the adjoints were defined as in (3) by the gradient of  $f$ , the direct differentiation was performed simultaneously with respect to all components of the discretized  $u$  so that the quantity  $\hat{u}$  occurring in (11) and (12) was in fact the identity matrix of order 12. Consequently,  $\hat{y}_k$  and  $\hat{y}_k$  had also 12 times as many components as the underlying  $y_k$  and  $\bar{y}_k$ , which are of the same size.

## 6 Summary, Conclusion and Outlook

We studied the convergence behavior of fixed point iterations for derivatives of implicit functions. These recurrences are generated in a completely mechanical fashion from a user supplied contractive fixed point solver for evaluating the implicit function. While the contractivity and thus the asymptotic convergence rate is inherited by the derived solvers there is a certain time lag. This is not really surprising since the equations for the adjoints  $\bar{y}$  and those for the feasible directions  $\hat{y}$  are dependent on  $y$  and both in turn impact the second order adjoint equation for  $\hat{y}$ . Mathematically we obtain Jordan blocks of size 2 for the double eigenvalues of the first derivative systems and of size 3 for the quadruple eigenvalues of the second order adjoint system. One does not obtain blocks of size 4 since the (3,2) sub-block in the big Jacobian system vanishes identically. Otherwise it would connect the two first derivative systems.

Generally if one were to iteratively evaluate derivatives of order  $d$  one can expect that the

relative errors compared to those of the underlying function iteration grows like  $k^d$ , where  $k$  is the iteration counter. In the context of constrained optimization one can expect that the correct values of reduced gradients (4) and Hessians (13) are obtained slower than feasibility so that optimality will be arrived at in the tangential fashion that is familiar from SQP calculations [8, 9, 10]. In fact when the state equation only be solved by a slowly convergent fixed point solver as we have assumed throughout it makes little sense to apply an SQP type algorithms. Instead one will prefer a so-called one-shot optimization strategy [11], where feasibility and optimality is achieved at the same time. We are currently investigating a piggy-back optimization scheme, where a third iteration updating the design variables  $u$  on the basis of approximate reduced gradient information is appended to (3).

## References

- [1] A. Griewank, C.H. Bischof, G.F. Corliss, A. Corle, and K. Williamson. Derivative Convergence of Iterative Equation Solvers. *Optimization Methods and Software*, 2:321-355, 1993.
- [2] A. Griewank and C. Fauré. Reduced functions, gradients and Hessians from fixed point iteration for state equations. *Numerical Algorithms*, 30(2):113-139, 2002.
- [3] Andreas Griewank, David Juedes, and Jean Utke. ADOL-C, A Package for the Automatic Differentiation of Algorithms Written in C/C++. *TOMS*, 22(2):131-167, 1996.
- [4] M. Hinze and T. Slawig. Adjoint gradients compared to gradients from algorithmic differentiation in instantaneous control of the navier-stokes equations. *Optimization Methods and Software*, 18(3):299-315, 2003.
- [5] A. Jameson. Optimum aerodynamic design using cfd and control theory. In *12th AIAA Computational Fluid Dynamics Conference, AIAA Paper 95-1729*, San Diego, CA, 1995, American Institute of Aeronautics and Astronautics.
- [6] B. Mohammadi and O. Pironneau. *Applied Shape Optimization for Fluids*. Numerical Mathematics and Scientific Computation. Clarendon Press, Oxford, 2001.
- [7] P.A. Newman, G.J.W. Hou, H.E. Jones, A.C. Taylor, and V.M. Korivi. Observations on computational methodologies for use in large-scale, gradient-based, multidisciplinary design incorporating advanced CFD codes. Technical Memorandum 104206. NASA Langley Research Center, February 1992. AVSCOM Technical Report 92-B-007.
- [8] J. Nocedal and S.J. Wright. *Numerical Optimization*, Springer Series in Operation Research. Springer Verlag, New York,....Tokyo, 1999.
- [9] E. W. Sachs. Control applications of reduced SQP methods. In R. Bulirsch and D. Kraft, editors, *Computational Optimal Control*, volume 115 of *Int. Series Num. Math.*, pages 89-104, Birkhäuser, 1994.
- [10] V. H. Schulz. Solving discretized optimization problems by partially reduced SQP methods. *Computing and Visualization in Science*, 1:83-96, 1998.
- [11] S. Tu'nessu, G. Kuruvila, and M.D. Salas. Aerodynamic design and optimization in one shot. In *30th AIAA Aerospace Sciences Meeting and Exhibit, AIAA Paper 91-0025*, Reno, Nevada, 1992. American Institute of Aeronautics and Astronautics.