

# 17 SOLUTION OF NONLINEAR EQUILIBRIUM EQUATIONS

We have seen that the finite element method for general nonlinear problems leads to the solution of nonlinear equations. The present chapter is therefore devoted to a discussion of various methods to solve nonlinear equations. Many of these methods have their origin not only in solid and structural mechanics, but also in nonlinear optimization theory and as general references, the reader may consult Bathe (1996), Belytschko *et al.* (2000), Crisfield (1991), Fletcher (1980), Luenberger (1984), Papadrakakis (1993) and Zienkiewicz and Taylor (1991) for relevant information.

The nonlinear equations of interest here are the equilibrium equations given by (16.10) and (16.11), i.e.

$$\boxed{\psi = 0} \quad (17.1)$$

where

$$\boxed{\psi = f_{int} - f} \quad (17.2)$$

and  $f$  denotes the *external forces*, i.e. the load on the body, defined by

$$f = \int_S N^T t dS + \int_V N^T b dV$$

whereas the *internal forces*  $f_{int}$  are defined by

$$f_{int} = \int_V B^T \sigma dV \quad (17.3)$$

Expression (17.1) holds for any body in equilibrium and our problem is to satisfy this equation.

In order to solve the boundary value problem in question, we have to consider the response of the actual material. Here, we will for illustration purposes assume elasto-plasticity i.e.

$$\boxed{\dot{\sigma} = D_t \dot{\epsilon}} \quad (17.4)$$

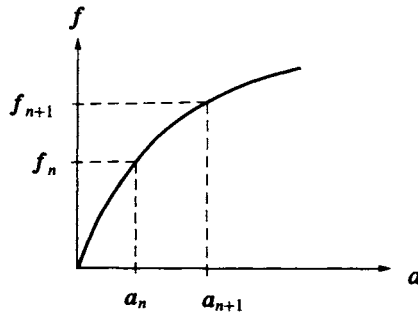


Figure 17.1: State  $n$  is known; we want to determinate state  $n + 1$ .

where  $D_t$  is the tangential stiffness equal to the elastic stiffness if elastic loading occurs and equal to the elasto-plastic tangential stiffness if plastic loading is present. Apart from this constitutive relation that is specific for elasto-plastic problems, all aspects and results that are related to the solution of the nonlinear equilibrium equations are general and they therefore also hold for viscoplasticity and creep, for instance. The equations of equilibrium are also referred to as the *global equations*, as they hold for the entire body, whereas the constitutive relations (17.4) are referred to as the *local equations*, as they hold for each material point.

The external loading  $f$  is assumed to be known and both (17.1) and (17.4) are nonlinear equation systems. However, since the elasto-plastic response of the material depends on the entire load history, we cannot simply solve (17.1) and (17.4) by imposing the entire external load directly. Instead, we have to adopt a *stepwise*, i.e. an *incremental solution procedure*, where the external loading is increased in small steps.

With reference to Fig. 17.1, we assume that we have reached state  $n$  where everything is known. That means the nodal displacements  $a_n$ , the external forces  $f_n$ , the strains  $\epsilon_n$  and the stresses  $\sigma_n$  are all known quantities. To start the process, the state  $n$  may be taken to be the state where the body is completely unloaded. The external load is now increased to  $f_{n+1}$  and we want to determine  $a_{n+1}$ ,  $\epsilon_{n+1}$  and  $\sigma_{n+1}$  at state  $n + 1$ .

Before a general approach to the solution of (17.1) is presented, it is of interest to discuss the simplest possible approach, the *Euler forward scheme*.

## 17.1 Euler forward scheme

In order to illustrate some important issues, we will start with the simple Euler forward scheme. Historically, this method was the first to be used for the solution of nonlinear finite element problems within solid mechanics. However,

it will turn out that this method possesses a fundamental drawback, but the illustration of this drawback may serve as a prelude to a more formal and correct manner of approaching our problem.

Since the constitutive relation (17.4) is given in an incremental fashion, it is tempting to differentiate (17.1) with respect to time to obtain

$$\int_V \mathbf{B}^T \dot{\boldsymbol{\sigma}} dV = \dot{\mathbf{f}} \quad (17.5)$$

where

$$\dot{\mathbf{f}} = \int_S \mathbf{N}^T \dot{\mathbf{t}} dS + \int_V \mathbf{N}^T \dot{\mathbf{b}} dV \quad (17.6)$$

According to (17.4) and our finite element approximation, we have

$$\dot{\boldsymbol{\sigma}} = \mathbf{D}_t \dot{\boldsymbol{\epsilon}} = \mathbf{D}_t \mathbf{B} \dot{\mathbf{a}}$$

Insertion of this expression into (17.5) provides

$$\mathbf{K}_t \dot{\mathbf{a}} = \dot{\mathbf{f}} \quad (17.7)$$

where the *tangential stiffness matrix*  $\mathbf{K}_t$  is defined by

$$\boxed{\mathbf{K}_t = \int_V \mathbf{B}^T \mathbf{D}_t \mathbf{B} dV} \quad (17.8)$$

It is evident that  $\mathbf{K}_t$  represents the current, i.e. the tangential stiffness of the body.

Referring to Fig. 17.1, everything is known at state  $n$ . The external load is increased to the known quantity  $\mathbf{f}_{n+1}$  and we want to determine  $\mathbf{a}_{n+1}$ ,  $\boldsymbol{\epsilon}_{n+1}$  and  $\boldsymbol{\sigma}_{n+1}$  at state  $n+1$ . Therefore, multiplying (17.7) by  $dt$  and integrating from state  $n$  to state  $n+1$ , we obtain

$$\int_{\mathbf{a}_n}^{\mathbf{a}_{n+1}} \mathbf{K}_t d\mathbf{a} = \mathbf{f}_{n+1} - \mathbf{f}_n \quad (17.9)$$

The fundamental problem that we are faced with is that  $\mathbf{K}_t$  is not a constant matrix; in fact,  $\mathbf{K}_t$  depends on the displacements, cf. (17.8), since  $\mathbf{D}_t$  does so. In addition, it is not even known how  $\mathbf{K}_t$  varies from state  $n$  to state  $n+1$ . However, at state  $n$  everything is known, i.e. also the constitutive matrix  $\mathbf{D}_t$  is known and the simplest thing we can do is to approximate  $\mathbf{K}_t$  in (17.9) by  $\mathbf{K}_t$  evaluated at state  $n$ . With the obvious notation

$$(\mathbf{K}_t)_n = \int_V \mathbf{B}^T (\mathbf{D}_t)_n \mathbf{B} dV$$

we then obtain approximately from (17.9) that

$$\boxed{(\mathbf{K}_t)_n (\mathbf{a}_{n+1} - \mathbf{a}_n) = \mathbf{f}_{n+1} - \mathbf{f}_n \quad \text{Euler forward}} \quad (17.10)$$

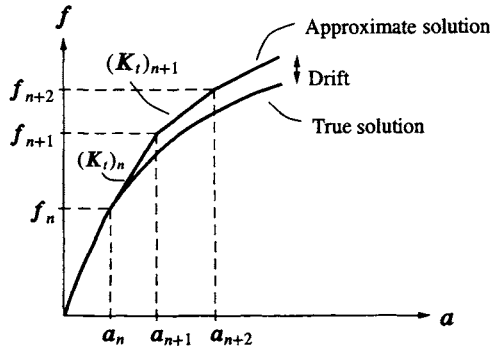


Figure 17.2: Euler forward scheme.

Except for  $a_{n+1}$ , everything is known and we can therefore solve this linear equation system to provide  $a_{n+1}$ . The approach obtained is termed the *Euler forward scheme*. The word 'forward' refers to the fact that we use our information about state  $n$  to determine state  $n+1$  by a direct extrapolation.

When  $a_{n+1}$  has been determined from (17.10), the corresponding strain is given by  $\epsilon_{n+1} = B a_{n+1}$ . Moreover, the stress state  $\sigma_{n+1}$  is obtained by integration of the constitutive relation (17.4) from state  $n$  to state  $n+1$ . In a symbolic manner, we write

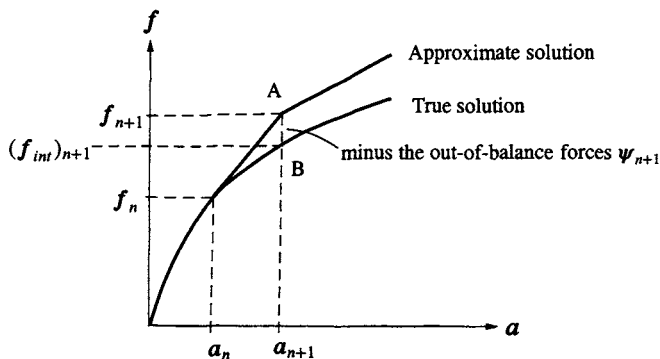
$$\sigma_{n+1} = \sigma_n + \int_{\epsilon_n}^{\epsilon_{n+1}} D_t d\epsilon \quad (17.11)$$

For points that respond elastically this integration is trivial. However, for points that respond in an elasto-plastic manner, the integration is far from being trivial. We will deal with this problem in detail in the next chapter; at the present stage, we simply accept that this integration can be performed.

Since (17.10) is certainly an approximate solution to our original problem, the repeated use of (17.10) to trace the response of the body is bound to introduce some errors. An illustration of this effect is shown in Fig. 17.2.

In this figure, we assume that the exact solution is known at state  $n$ . Increasing the external load from  $f_n$  to  $f_{n+1}$ , (17.10) then determines the nodal displacements  $a_{n+1}$ . Since  $(K_t)_n$  is the tangential stiffness matrix at state  $n$ , we obtain the over-shooting effect shown in Fig. 17.2. If this process is now repeated with the state  $n+1$  considered to be the known state, we obtain the response to the next step also shown in Fig. 17.2. It is apparent that use of the Euler forward scheme gives rise to a certain drift from the true solution.

This immediately raises the question whether it is possible to correct this drift and thereby improve our solution. In order to do so, we must find a quantitative expression for the drift and this problem is complicated by the fact that, in general, we do not know the true solution. The key point in this discussion



**Figure 17.3:** Discussion of drift and equilibrium;  $(f_{int})_{n+1}$  denotes the internal forces corresponding to the stresses  $\sigma_{n+1}$ . These stresses are calculated from the nodal displacements  $a_{n+1}$ .

turns out to be that of equilibrium.

Referring to (17.2), it is recalled that  $f_{int}$  represents the internal forces that the stresses give rise to. Since  $f$  denotes the known external forces, (17.1) states that equilibrium of the body requires that the external forces be equal to the internal forces. With the stresses  $\sigma_{n+1}$  determined by (17.11), we can calculate the internal forces from (17.3) according to

$$(f_{int})_{n+1} = \int_V B^T \sigma_{n+1} dV$$

The internal forces  $(f_{int})_{n+1}$  correspond to the stresses  $\sigma_{n+1}$  and these stresses are calculated from the nodal displacements  $a_{n+1}$ . Since the integration (17.11) may for the time being be assumed to be exact, the internal forces  $(f_{int})_{n+1}$  take the value illustrated in Fig. 17.3. It is evident that  $(f_{int})_{n+1} - f_{n+1}$ , i.e. the drift indicated in Fig. 17.2, is expressed by the fact that the Euler forward scheme does not fulfill equilibrium of the body. The drift  $(f_{int})_{n+1} - f_{n+1}$  is also called the *out-of-balance* or *residual forces* and the equilibrium condition (17.1) states that the out-of-balance forces, i.e. the residual forces, must be zero.

We have seen that the Euler forward scheme results in a certain drift from the true solution and this drift manifests itself in the form that the residual forces are different from zero, i.e. the equilibrium conditions for the body are not fulfilled.

Against this background, it seems natural to evaluate the nonlinear equilibrium equations (17.1) in more detail and investigate how we can ensure that our solution fulfills these conditions. Before we turn to this subject, we summarize the Euler forward algorithm as shown in Box 17.1. In this box, the means to enforce the boundary conditions has not been indicated; this aspect is dealt with in Section 17.6.

---

**Box 17.1** Euler forward algorithm
 

---

- *Initiation of quantities*  
 $\mathbf{a}_0 = \mathbf{0}; \quad \boldsymbol{\varepsilon}_0 = \mathbf{0} \quad \boldsymbol{\sigma}_0 = \mathbf{0}; \quad f_0 = 0$
  - *For load step  $n = 0, 1, 2, \dots, N_{max}$* 
    - *Determine new load level  $f_{n+1}$*
    - *Calculate  $\mathbf{K}_t = \int_V \mathbf{B}^T (\mathbf{D}_t)_n \mathbf{B} dV$*
    - *Calculate  $\mathbf{a}_{n+1}$  from  $\mathbf{K}_t (\mathbf{a}_{n+1} - \mathbf{a}_n) = f_{n+1} - f_n$*
    - *Calculate  $\boldsymbol{\varepsilon}_{n+1} := \mathbf{B} \mathbf{a}_{n+1}$*
    - *Determine  $\boldsymbol{\sigma}_{n+1}$  by integration of the constitutive equations (see next chapter)*
    - *Accept quantities*  
 $\mathbf{a}_{n+1}; \quad \boldsymbol{\varepsilon}_{n+1}; \quad \boldsymbol{\sigma}_{n+1}$
  - *End load step loop*
- 

## 17.2 General iteration format

We have argued for the requirement that our solution fulfills the equilibrium equations (17.1). Since these equations are nonlinear, this section is devoted to a discussion of some general aspects relating to the solution of nonlinear equations and in the next sections, we will then adopt these viewpoints to the nonlinear FE equations.

Suppose that we have some unknowns collected in the column matrix  $\mathbf{a}$ . These unknowns are determined by the following nonlinear equation system

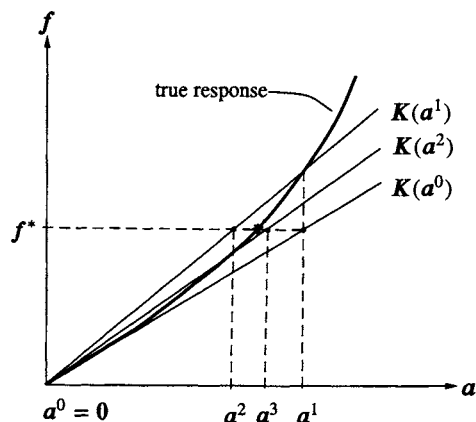
$$\mathbf{a} = \mathbf{F}(\mathbf{a}) \quad (17.12)$$

To solve this equation system, a common *iteration scheme* is

$$\mathbf{a}^i = \mathbf{F}(\mathbf{a}^{i-1}); \quad i = 1, 2, \dots \quad (17.13)$$

where  $i$  denotes the number of iterations and where  $\mathbf{a}^0$  is the so-called *start vector* that must be specified by us ( $\mathbf{a}^{i-1} = \mathbf{a}^0$  for  $i = 1$ ).

Let the true solution of (17.12) be denoted by  $\mathbf{z}$  and the question then arises whether the iteration scheme (17.13) converges towards the true solution  $\mathbf{z}$ . A strict mathematical proof that defines the conditions for which the iteration scheme (17.13) converges towards the true solution, is given, for instance, by Dahlquist and Björk (1974). However, it turns out to be difficult to make practical use of these mathematical conditions and we will therefore simply state the



**Figure 17.4:** Convergent iteration scheme. True solution for external force  $f^*$  is indicated by (\*).

*convergence theorem* in the following verbal form

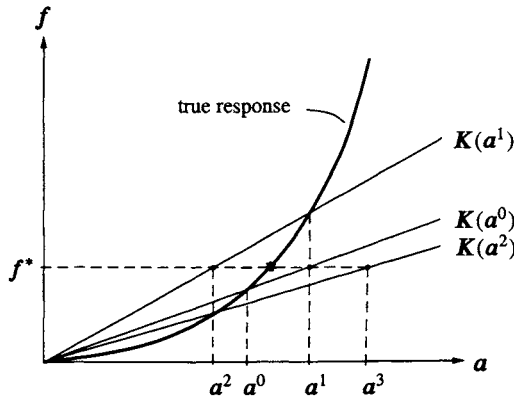
*If the start vector  $a^0$  is sufficiently close to the true solution of (17.12), then the iteration scheme (17.13) will converge towards this true solution*

What in this formulation is meant by 'sufficiently close' cannot be given an explicit formulation that is applicable in practice. However, it emphasizes an important point that also arises in nonlinear finite element calculations, namely that for a given elasto-plastic problem, say, we cannot, in general, be sure that a given solution scheme will provide the solution. Instead, we may have to modify the load steps, the number of iterations and even be forced to adopt another solution scheme in order to obtain a solution. In general therefore, nonlinear finite elements calculations are far from being trivial, but a solid knowledge of the underlying theory and experience in the solution of similar problems greatly enhance the possibility of achieving a solution.

To illustrate the procedure, we will consider the following equation system

$$K(a)a = f$$

where  $K$  is a nonlinear global stiffness matrix,  $f$  the external load and  $a$  the nodal displacements; the stiffness matrix  $K(a)$  can be viewed as the secant stiffness. This expression is typical for elastic problems involving geometrical nonlinearities and in this case one may often have the response shown in Fig. 17.4. Here, the structure stiffens with increasing loading and an example is the response of a laterally loaded plate with fixed supports, which stiffens due to the membrane forces created by large deformations. According to (17.12)



**Figure 17.5:** Divergent iteration scheme. True solution for external force  $f^*$  is indicated by (\*).

and (17.13), we adopt the following iteration scheme after the external load has been increased from zero to  $f = f^*$

$$a^i = [K(a^{i-1})]^{-1} f^* ; \quad i = 1, 2, \dots$$

Choosing the start vector as  $a^0 = 0$ , then  $K(a^0) = K(0)$  becomes equal to the initial elastic stiffness, and we then obtain the convergent iteration scheme shown in Fig. 17.4.

To illustrate that even small changes in the response may create a divergent iteration scheme, consider the response shown in Fig. 17.5. This response is quite close to that shown in Fig. 17.4, except that the stiffening effect is larger. The starting vector  $a^0$  is now taken to be different from zero; in fact the start vector is rather close to the exact solution. Nevertheless, Fig. 17.5 shows a divergent iteration scheme and it is somewhat surprising that the qualitatively similar responses shown in Figs. 17.4 and 17.5 can give rise to completely different iteration courses.

### 17.3 Standard iteration format for equilibrium iterations - Iteration matrix

Having discussed the general iteration format which started with the nonlinear equation system (17.12) and resulted in the iteration scheme (17.13), we will now see how the nonlinear equilibrium equations arising in nonlinear finite element calculations can be treated in the same manner.

We have seen that the backbone in nonlinear FE analysis is the fulfillment of the nonlinear equilibrium equations (17.1). Since the stresses depend on the



nodal displacements  $\mathbf{a}$ , and since we want to fulfill equilibrium at a given fixed external loading, we have from (17.1)-(17.3)

$$\boldsymbol{\psi}(\mathbf{a}) = \mathbf{0} \quad (17.14)$$

This relation may also be expressed in the form of the homogeneous equation system

$$\mathbf{0} = -(\mathbf{A}(\mathbf{a}))^{-1}\boldsymbol{\psi}(\mathbf{a}) \quad (17.15)$$

where  $\mathbf{A}^{-1}$  is a square matrix. In the most general case,  $\mathbf{A}^{-1}$  may also depend on the unknown  $\mathbf{a}$ -values. In order that (17.15) should provide the trivial solution  $\boldsymbol{\psi} = \mathbf{0}$  only, we must require

$$\boxed{\det \mathbf{A}^{-1} \neq 0} \quad (17.16)$$

This requirement is certainly fulfilled if  $\mathbf{A}^{-1}$  is positive definite, but even if  $\mathbf{A}^{-1}$  is not positive definite, its use in (17.15) is acceptable as long as (17.16) holds.

Expression (17.15) may even be formulated as

$$\mathbf{a} = \mathbf{a} - (\mathbf{A}(\mathbf{a}))^{-1}\boldsymbol{\psi}(\mathbf{a})$$

or

$$\mathbf{a} = \mathbf{F}(\mathbf{a}) \quad \text{where} \quad \mathbf{F}(\mathbf{a}) = \mathbf{a} - (\mathbf{A}(\mathbf{a}))^{-1}\boldsymbol{\psi}(\mathbf{a}) \quad (17.17)$$

By these manipulations, we have retrieved formulation (17.12).

In accordance with the general iteration scheme (17.13), we then obtain from (17.17) that

$$\mathbf{a}^i = \mathbf{a}^{i-1} - (\mathbf{A}(\mathbf{a}^{i-1}))^{-1}\boldsymbol{\psi}(\mathbf{a}^{i-1})$$

which may be written as

*Standard iteration format*

$$\mathbf{A}(\mathbf{a}^{i-1})(\mathbf{a}^i - \mathbf{a}^{i-1}) = -\boldsymbol{\psi}(\mathbf{a}^{i-1}); \quad i = 1, 2, \dots \quad (17.18)$$

*where  $\mathbf{A}$  is the iteration matrix*

All quantities denoted by  $i - 1$  are known quantities and (17.18) therefore determines the new nodal displacement estimate  $\mathbf{a}_i$ . This iteration scheme is the vehicle by which we solve (17.14) and it will be referred to as *the standard iteration format*; the matrix  $\mathbf{A}$  is called the *iteration matrix*.

The iteration matrix  $\mathbf{A}$  must be chosen by us and we will later see how the format (17.18) enables us to derive a very large group of iteration methods applicable to nonlinear FE calculations in a very simple and convenient manner. Before we probe further into this topic, it is worthwhile to scrutinize the restrictions that we must place on the iteration matrix  $\mathbf{A}$  in more detail. We have

already indicated restriction (17.16). Suppose that equilibrium is satisfied, i.e.  $\psi(\mathbf{a}^{i-1}) = \mathbf{0}$ . Then the correct solution has been achieved and the iteration scheme (17.18) should then imply no further changes of the nodal displacements, i.e.  $\mathbf{a}^i = \mathbf{a}^{i-1}$ . However, the only possibility that  $\mathbf{A}(\mathbf{a}^i - \mathbf{a}^{i-1}) = \mathbf{0}$  provides this solution is

$$\boxed{\det \mathbf{A} \neq 0} \quad (17.19)$$

Apart from the restrictions given by (17.16) and (17.19), the iteration matrix may be chosen arbitrarily, but we will experience that especially restriction (17.19) has important consequences for nonlinear finite element calculations.

Referring to our fundamental problem illustrated in Fig. 17.1, we start from state  $n$  where equilibrium is fulfilled and where the nodal displacements  $\mathbf{a}_n$ , the strains  $\boldsymbol{\epsilon}_n$ , the stresses  $\boldsymbol{\sigma}_n$  and the external loading  $\mathbf{f}_n$  are all known. The external loading is then increased to  $\mathbf{f}_{n+1}$  and we then want to determine the corresponding displacements  $\mathbf{a}_{n+1}$ , the strains  $\boldsymbol{\epsilon}_{n+1}$  and thereby also the stresses  $\boldsymbol{\sigma}_{n+1}$  (obtained by an integration of the constitutive equations). The purpose of the iteration scheme (17.18) is to fulfill equilibrium at state  $n + 1$ . Since the external loading is given by  $\mathbf{f}_{n+1}$ , the out-of-balance forces  $\psi(\mathbf{a}^{i-1})$  defined by (17.2)-(17.3) become

$$\boxed{\psi(\mathbf{a}^{i-1}) = \int_V \mathbf{B}^T \boldsymbol{\sigma}^{i-1} dV - \mathbf{f}_{n+1}} \quad (17.20)$$

Whereas the external loading  $\mathbf{f}_{n+1}$  is specified by us, the stress  $\boldsymbol{\sigma}^{i-1}$  must be obtained by integration of the constitutive relations (17.4). The last state where the stresses were known and where the equilibrium equations were satisfied was state  $n$ . For each point in the body,  $\boldsymbol{\sigma}^{i-1}$  is therefore obtained by the following symbolic integration

$$\begin{aligned} \boldsymbol{\sigma}^{i-1} &= \boldsymbol{\sigma}_n + \mathbf{D}(\boldsymbol{\epsilon}^{i-1} - \boldsymbol{\epsilon}_n) \quad \text{if the point behaves elastically} \\ \boldsymbol{\sigma}^{i-1} &= \boldsymbol{\sigma}_n + \int_{\boldsymbol{\epsilon}_n}^{\boldsymbol{\epsilon}^{i-1}} \mathbf{D}^p d\boldsymbol{\epsilon} \quad \text{if the point behaves plastically} \end{aligned} \quad (17.21)$$

The specific manner by which the integration is performed when plastic behavior occurs is dealt with in the next chapter; at the present stage, we simply accept that this integration can be performed. In practice, the determination of the stresses  $\boldsymbol{\sigma}^{i-1}$  by means of (17.21) is not performed for all points in the body, but only for the *Gauss points*.

To start the general iteration scheme (17.18) for  $i = 1$ , we need to specify a starting value of  $\mathbf{a}$ , i.e. we have to specify  $\mathbf{a}^0$ . The most recent known value of  $\mathbf{a}$  is  $\mathbf{a}_n$  and it is therefore appropriate to choose  $\mathbf{a}^0 = \mathbf{a}_n$  and thereby  $\boldsymbol{\epsilon}^0 = \boldsymbol{\epsilon}_n$  which, according to (17.21) implies  $\boldsymbol{\sigma}^0 = \boldsymbol{\sigma}_n$ . Likewise, as the iteration matrix

$\mathbf{A}^0$  we take the one given by  $\mathbf{A}_n$ . The starting conditions for the iteration scheme therefore become

$$\mathbf{a}^0 = \mathbf{a}_n ; \quad \boldsymbol{\sigma}^0 = \boldsymbol{\sigma}_n ; \quad \mathbf{A}^0 = \mathbf{A}_n \quad (17.22)$$

From (17.20), we then obtain

$$\boldsymbol{\psi}^0 = \int_V \mathbf{B}^T \boldsymbol{\sigma}_n dV - \mathbf{f}_{n+1} \quad (17.23)$$

By means of (17.18), we have within each load step created a series of iterations which gradually improve the solution. The quantity that controls the iterations is the out-of-balance forces  $\boldsymbol{\psi}(\mathbf{a}^{i-1})$ , which measure the difference between the internal forces  $\int_V \mathbf{B}^T \boldsymbol{\sigma}^{i-1} dV$  and the external forces  $\mathbf{f}_{n+1}$ . When the out-of-balance forces approach zero, the correction to the nodal displacements also approaches zero. In practice therefore, the iteration scheme (17.18) is stopped when the out-of-balance forces become smaller than a certain amount specified by the user of the FE program. We shall return to this subject in Section 17.7. When  $\boldsymbol{\psi}(\mathbf{a}^{i-1})$  is small enough, we accept the solution  $\mathbf{a}^{i-1}$ , i.e

<p style="margin: 0;">When convergence is accepted</p> $\mathbf{a}_{n+1} = \mathbf{a}^{i-1} ; \quad \boldsymbol{\epsilon}_{n+1} = \boldsymbol{\epsilon}^{i-1} ; \quad \boldsymbol{\sigma}_{n+1} = \boldsymbol{\sigma}^{i-1}$	(17.24)
--	---------

Since the objective of the iterations indicated in (17.18) is to fulfill the equilibrium equations, these iterations are called *equilibrium iterations*. Moreover, since the external load is increased incrementally and as iterations are performed within each load step, this approach is often called an *incremental-iterative approach*.

Let us return to the integration of the constitutive equations given by (17.21). The integration limits turn out to be very important and we observe that they are given by  $\boldsymbol{\epsilon}_n$ , which is the last state where the equilibrium conditions were satisfied, and  $\boldsymbol{\epsilon}^{i-1}$  which is the current value of the strains. That is, with the formulation (17.21), we deliberately ignore the intermediate states between  $\boldsymbol{\epsilon}_n$  and  $\boldsymbol{\epsilon}^{i-1}$  since these intermediate states may, for instance, involve false elastic unloading excursions. If we performed the integration of the constitutive equations by first integrating from  $\boldsymbol{\epsilon}_n$  to  $\boldsymbol{\epsilon}^1$ , then from  $\boldsymbol{\epsilon}^1$  to  $\boldsymbol{\epsilon}^2$  and so on, we would force the stress state  $\boldsymbol{\sigma}^{i-1}$  to depend on some possibly false intermediate states, i.e. false load history and this pitfall is avoided by formulation (17.21). Apart from that in the next chapter, we shall see how integration of (17.21) is performed in practice.

It may come as some surprise that the only point where the constitutive equations are involved is in the integration scheme (17.21). Moreover, the iteration strategy (17.18) has the very important property that possible non-zero out-of-balance forces in one load step are automatically transferred to the next load step. That is, there is no possibility for an accumulation of non-zero out-of-balance forces.

It is of significant interest that all iteration strategies proposed in the literature and used to solve nonlinear FE equations are, in principle, embraced by the scheme (17.18) and (17.20)-(17.24). The key point is the selection of the iteration matrix  $\mathbf{A}$  where

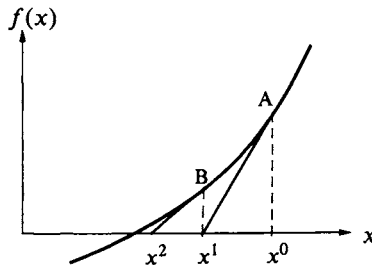
*The choice of a specific iteration matrix  $\mathbf{A}$  represents the establishment of a specific solution strategy*

(17.25)

We recall that the only restrictions on the iteration matrix  $\mathbf{A}$  are given by (17.16) and (17.19) and if (17.18) is to be an expression that is dimensionally correct,  $\mathbf{A}$  must be a stiffness matrix; except for that, each choice of  $\mathbf{A}$  represents a valid solution strategy. Certainly, this is not to say that all choices of  $\mathbf{A}$  are equally promising, but the statement (17.25) is the vehicle that produces an arsenal of different methods where specific advantages and shortcomings are related to each choice of the iteration matrix  $\mathbf{A}$ . With this general framework, we will now discuss some classical and some more recent iteration strategies.

## 17.4 Newton-Raphson scheme

It was mentioned that the choice of different iteration matrices  $\mathbf{A}$  in the standard iteration scheme (17.18) creates different solution strategies. The choice of  $\mathbf{A}$  is subject to the constraints (17.16) and (17.19) and that it should represent some stiffness matrix, but apart from that we can make any choice. Evidently, not all choices create an efficient solution scheme and to obtain some kind of feeling for what is an appropriate choice, we will now identify the iteration matrix  $\mathbf{A}$  that corresponds to the well-known *Newton-Raphson scheme*. Often this method is simply called *Newton's method*, but it was derived simultaneously by Raphson; Bičanić and Johnson (1979) give the relevant historical background.



**Figure 17.6:** Newton-Raphson strategy for a one-dimensional problem.

For a one-dimensional problem, the essence of the Newton-Raphson strategy is illustrated in Fig. 17.6. The problem is to identify the solution to the

nonlinear equation  $f(x) = 0$ . A starting value  $x^0$  is guessed by us. At the corresponding point A on the curve  $f(x)$ , the tangent is determined and this tangent is extrapolated to obtain the next estimate  $x^1$  for the solution. This process is then repeated so that the tangent at point B provides the next estimate  $x^2$ , etc.

It appears that the fundamental idea of the Newton-Raphson approach is to linearize the nonlinear function about a given point. This means that the nonlinear function is approximated by a Taylor expansion about the point in question and in this Taylor expansion, terms higher than the linear ones are ignored.

In our case, the nonlinear multi-dimensional function is given by the equilibrium equations, i.e.

$$\psi(a) = 0$$

where

$$\psi(a) = \int_V B^T \sigma dV - f \quad (17.26)$$

and the external forces  $f$  are known and fixed whereas the stresses  $\sigma$  depend on the nodal displacements  $a$ . Assume now that the approximation  $a^{i-1}$  to the true solution  $a$  has been established. Ignoring higher-order terms, a Taylor expansion of  $\psi$  about  $a^{i-1}$  yields

$$\psi(a^i) = \psi(a^{i-1}) + \left(\frac{\partial \psi}{\partial a}\right)^{i-1} (a^i - a^{i-1}) \quad (17.27)$$

This expression provides the linearized approximation to the true expression for  $\psi(a^i)$ ; it therefore represents the tangent to the curve at point  $a^{i-1}$ . Similar to the one-dimensional case, we require  $\psi(a^i) = 0$  and it then follows from (17.27) that

$$0 = \psi(a^{i-1}) + \left(\frac{\partial \psi}{\partial a}\right)^{i-1} (a^i - a^{i-1}) \quad (17.28)$$

To proceed further, we have to identify the derivative  $\partial \psi / \partial a$ . Since the external loading is fixed, (17.26) implies

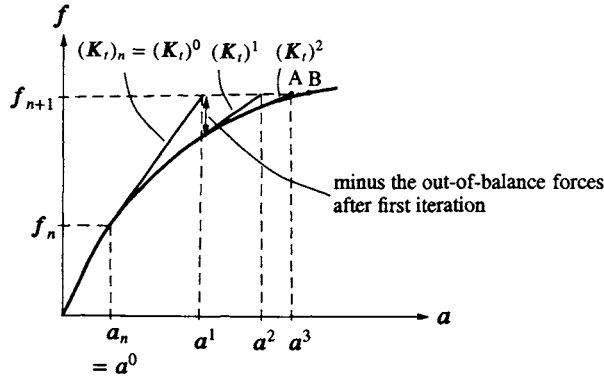
$$\frac{\partial \psi}{\partial a} = \int_V B^T \frac{d\sigma}{da} dV \quad (17.29)$$

From the constitutive relation  $\dot{\sigma} = D_t \dot{\epsilon}$ , cf. (17.4), and noting that the variable now is  $a$ , we obtain

$$d\sigma = D_t d\epsilon = D_t B da$$

i.e.

$$\frac{d\sigma}{da} = D_t B \quad (17.30)$$



**Figure 17.7:** Newton-Raphson scheme showing the equilibrium iterations. Point B is the true solution and point A is the solution obtained after three iterations.

Insertion of (17.30) into (17.29) yields

$$\frac{\partial \psi}{\partial a} = K_t \quad \text{where} \quad K_t = \int_V B^T D_t B dV \quad (17.31)$$

where  $K_t$  is the tangent stiffness matrix of the body, cf. (17.8).

With this result, (17.28) takes the form

$$(K_t)^{i-1}(a^i - a^{i-1}) = -\psi(a^{i-1}) \quad (17.32)$$

and a comparison with the standard iteration scheme (17.18) shows that

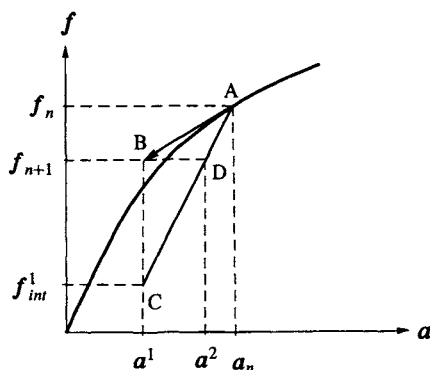
$$\boxed{A = K_t \Rightarrow \text{Newton-Raphson scheme}} \quad (17.33)$$

Therefore, this choice of the iteration matrix  $A$  results in the well-known Newton-Raphson approach.

For the first iteration  $i = 1$ , we find with (17.23) and the starting conditions (17.22) that

$$(K_t)_n(a^1 - a_n) = f_{n+1} - \int_V B^T \sigma_n dV \quad \text{first iteration}$$

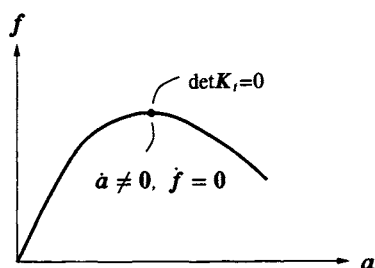
If equilibrium is fulfilled at state  $n$  then  $f_n = \int_V B^T \sigma_n dV$  and we immediately observe that the first iteration is, in fact, identical to the Euler forward scheme, cf. (17.10). From (17.32) follows that in each equilibrium iteration, the Newton-Raphson approach makes use of the current tangential stiffness matrix and this feature is illustrated in Fig. 17.7. The Newton-Raphson approach is one of the most often used solution schemes and the reason is the fast convergence as illustrated in Fig. 17.7. Since the Newton-Raphson method starts with a Euler forward prediction and then continues with successive corrections



**Figure 17.8:** Newton-Raphson scheme during unloading of elasto-plastic body. Correct solution is obtained after two iterations.

to fulfill equilibrium, the terminology of a *predictor-corrector scheme* is often used.

In Fig. 17.7, we assumed that plastic loading occurs and it may be of interest to investigate the Newton-Raphson approach when elastic unloading occurs and Fig. 17.8 illustrates this issue. It appears that the correct response will be predicted after two iterations. First, unloading occurs along AB using the tangential stiffness at point A. Observing that, in reality, elastic unloading occurs along AC, the internal forces corresponding to the displacements  $a^1$  are given by point C. That is, at that stage the out-of-balance forces  $\psi(a^1)$  become with (17.20)  $\psi(a^1) = f_{int}^1 - f_{n+1} = -|BC|$ . When these residual forces are applied in (17.32) in the next iteration where  $i = 2$  and observing that the current tangential stiffness now amounts to the elastic stiffness at point C, we arrive at the correct point D after the second iteration.



**Figure 17.9:** The Newton-Raphson scheme will have difficulties close to a peak load.

Since the tangential matrix  $K_t$  varies with the loading, it is of interest to investigate whether at some state in the Newton-Raphson scheme we may violate

---

**Box 17.2** Newton-Raphson algorithm
 

---

- *Initiation of quantities*

$$\mathbf{a}_0 = \mathbf{0}; \quad \boldsymbol{\varepsilon}_0 = \mathbf{0}; \quad \boldsymbol{\sigma}_0 = \mathbf{0}; \quad \mathbf{f}_0 = \mathbf{0}; \quad \mathbf{f}_{int} = \mathbf{0}$$
  - *For load step*  $n = 0, 1, 2, \dots, N_{max}$ 
    - *Determine new load level*  $\mathbf{f}_{n+1}$
    - *Initiation of iteration quantities*

$$\mathbf{a}^0 := \mathbf{a}_n$$
    - *Iterate*  $i = 1, 2, \dots$  *until*  $|\boldsymbol{\psi}|_{norm} = |\mathbf{f}_{int} - \mathbf{f}_{n+1}|_{norm} < \epsilon_{\text{epsilon}}$ 
      - *Calculate*  $\mathbf{K}_i = \int_V \mathbf{B}^T \mathbf{D}_i^{i-1} \mathbf{B} dV$
      - *Calculate*  $\mathbf{a}^i$  *from*  $\mathbf{K}_i(\mathbf{a}^i - \mathbf{a}^{i-1}) = \mathbf{f}_{n+1} - \mathbf{f}_{int}$
      - *Calculate*  $\boldsymbol{\varepsilon}^i := \mathbf{B}\mathbf{a}^i$
      - *Determine*  $\boldsymbol{\sigma}^i$  *by integration of the constitutive equations*  
(see next chapter)
      - *Calculate internal forces*  $\mathbf{f}_{int} = \int_V \mathbf{B}^T \boldsymbol{\sigma}^i dV$
    - *End iteration loop*
    - *Accept quantities*

$$\mathbf{a}_{n+1} := \mathbf{a}^i; \quad \boldsymbol{\varepsilon}_{n+1} := \boldsymbol{\varepsilon}^i; \quad \boldsymbol{\sigma}_{n+1} := \boldsymbol{\sigma}^i; \quad \mathbf{f}_{int}$$
  - *End load step loop*
- 

the important restriction given by (17.19), which with (17.33) becomes

$$\det \mathbf{K}_i \neq 0$$

This restriction is violated if  $\det \mathbf{K}_i = 0$  which implies that the homogeneous equation system  $\mathbf{K}_i \dot{\mathbf{a}} = \mathbf{0}$  possesses a non-trivial  $\dot{\mathbf{a}}$ -solution. As illustrated in Fig. 17.9, this situation corresponds to the peak load where we have  $\dot{\mathbf{a}} \neq \mathbf{0}$ ,  $\dot{\mathbf{f}} = \mathbf{0}$  which satisfy the equation  $\mathbf{K}_i \dot{\mathbf{a}} = \dot{\mathbf{f}}$  from (16.17) when  $\det \mathbf{K}_i = 0$ . Use of the Newton-Raphson scheme will therefore imply increasing difficulties when a peak load is approached. If the response of the body, as illustrated in Fig. 17.9, exhibits a softening branch, great difficulties are encountered when we try to trace the response over the peak load and into the softening region. This latter observation is relevant, not only for the Newton-Raphson scheme, but also for other methods and we shall return to this subject later on.

We finally observe that in every Newton-Raphson iteration, a new  $\mathbf{K}_i$ -matrix needs to be established and that, in principle, the inverse matrix  $(\mathbf{K}_i)^{-1}$  also



needs to be identified, cf. (17.32), i.e. each such iteration is costly. With this remark, we may summarize the properties of the Newton-Raphson scheme as follows:

- \* *Newton-Raphson works well both in loading and unloading*
  - \* *Newton-Raphson provides a fast convergence*
  - \* *problems may occur close to peak points*
  - \* *every Newton-Raphson iteration is costly*
- (17.34)

The Newton-Raphson algorithm is summarized in Box. 17.2; here the means to impose the boundary conditions have not been indicated, see Section 17.6.

## 17.5 Initial stiffness and modified Newton-Raphson schemes

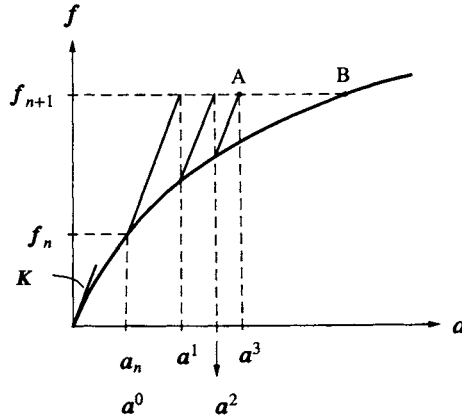
It was observed that the choice  $\mathbf{A} = \mathbf{K}$ , leads to the well-known and efficient Newton-Raphson scheme. However, we have previously stressed the freedom that we have when choosing the iteration matrix  $\mathbf{A}$  in the standard iteration scheme (17.18). Against this background, it is not surprising that a number of choices exists for the iteration matrix which are modifications of the Newton-Raphson approach. Since such modifications exist, the true Newton-Raphson method is occasionally referred to as the *full* Newton-Raphson scheme. Scrutinizing the standard iteration scheme (17.18), it follows that any  $\mathbf{A}$ -matrix must have the same dimension as the stiffness matrix of the body.

Recognizing that every Newton-Raphson iteration is costly, it is tempting to choose a constant  $\mathbf{A}$ -matrix in all iterations. The most evident choice would then be  $\mathbf{A} = \mathbf{K}$  where  $\mathbf{K}$  is the linear elastic stiffness matrix. Since  $\mathbf{K}$  is the initial stiffness of the body, this choice leads to the *initial stiffness method*, i.e.

$$\mathbf{A} = \mathbf{K} \Rightarrow \text{Initial stiffness scheme}$$

Since the elastic stiffness matrix  $\mathbf{K}$  is positive definite, we have  $\det \mathbf{K} \neq 0$  and none of the restrictions given by (17.16) and (17.19) are violated.

The performance of this approach is illustrated in Fig. 17.10 and a comparison with Fig. 17.7 clearly illustrates that the initial stiffness method simplifies the calculations considerably, but at the expense of a much slower convergence. This is especially true when the structure approaches its peak load, where the force-displacement curve is flat. For elastic unloading of the body, it is evident that the initial stiffness method will provide the correct response in just one



**Figure 17.10:** Initial stiffness method. Point B is the true solution and point A is the solution obtained after three iterations.

iteration. We are then led to the following conclusions

- \* *Initial stiffness approach works well both in loading and unloading*
- \* *every iteration is cheap*
- \* *Initial stiffness approach converges slowly*

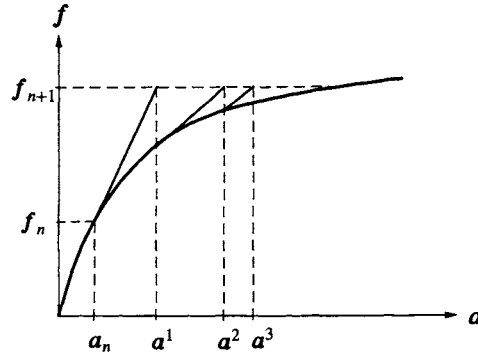
The initial stiffness method may be recast into other equivalent forms that have been proposed in the literature. With Hooke's law  $\sigma = D(\epsilon - \epsilon^p) = DBa - D\epsilon^p$ , the internal forces become

$$f_{int}^{i-1} = \int_V B^T \sigma^{i-1} dV = K a^{i-1} - \int_V B^T D(\epsilon^p)^{i-1} dV$$

where  $K = \int_V B^T D B dV$  is the elastic stiffness matrix. Insertion of the expression above into the general iteration scheme (17.18) and (17.20) with  $A = K$  gives

$$K a^i = f_{n+1} + \int_V B^T D(\epsilon^p)^{i-1} dV \quad (17.35)$$

This formulation was suggested by Argyris (1965b) under the name *initial strain method*. The reason for this terminology is that linear elasticity with initial strains is given by  $\sigma = D(\epsilon - \epsilon^o)$  where  $\epsilon^o$  are the initial strains, cf. (4.62). It appears that the plastic strains  $\epsilon^p$  may be interpreted as initial strains. On the other hand, if we define the quantity  $\sigma^o$  by  $\sigma^o = -D\epsilon^p$  then Hooke's law can be written as  $\sigma = D\epsilon + \sigma^o$ , i.e. for zero strains, we have  $\sigma = \sigma^o$  and  $\sigma^o$  may therefore be viewed as initial stresses. Replacing  $D\epsilon^p$  by  $-\sigma^o$ , the



**Figure 17.11:** Modified Newton-Raphson. Updating of the tangential stiffness matrix after the first iteration in each load step.

scheme (17.35) was proposed by Zienkiewicz *et al.* (1969) under the name: *initial stress method*. However, as we have seen, both the initial strain and the initial stress method are identical to the initial stiffness method; it is only the physical interpretation that differs.

A compromise between full Newton-Raphson, where updating of the tangential stiffness matrix  $\mathbf{K}_t$  occurs in every iteration, and the initial stiffness method, where no updating occurs at all, is the modification where updating is performed only once in each load step. This approach is called *modified Newton-Raphson*. Intuitively, it seems reasonable to perform this updating at the beginning of the load step. However, we experienced in Fig. 17.8 that full Newton-Raphson predicts the correct response during elastic unloading after two iterations. Therefore, the approach suggested above would cause problems during unloading. Instead, the updating of the tangential stiffness matrix is chosen to occur after the first equilibrium iteration in each load step. This approach is illustrated in Fig. 17.11.

Finally, we may mention the so-called *self-correcting procedure* of Stricklin *et al.* (1971) and further elaborated on by Stricklin and Haisler (1977). Essentially, it consists of using the full Newton-Raphson approach, but make just one equilibrium iteration in every load step. That is, the general iteration scheme (17.18) is used with  $i = 1$  and with  $\mathbf{A} = \mathbf{K}_t$ . We have previously emphasized that the scheme (17.18) has the advantage that possible non-zero out-of-balance forces in one load step are automatically transferred to the next load step and the self-correcting procedure takes full advantage of this property.

To illustrate this issue, assume that the numerical procedure has determined point A in Fig. 17.12 to be the response at state  $n$ . For the nodal displacements  $\mathbf{a}_n$ , the corresponding internal forces are given by  $(\mathbf{f}_{int})_n = \int_V \mathbf{B}^T \boldsymbol{\sigma}_n dV$ , i.e. point P. The out-of-balance forces are therefore given by  $(\mathbf{f}_{int})_n - \mathbf{f}_n = -|\mathbf{AP}|$ . When the external loading is increased to  $\mathbf{f}_{n+1}$ , we obtain from (17.18) with



As an example, we may refer to the Newton-Raphson scheme given by (17.32) where the tangential stiffness matrix  $\mathbf{K}_t$  defined by (17.31) enters the equation system. When establishing  $\mathbf{K}_t$  in this manner, no considerations were taken of the kinematic boundary conditions; therefore, (17.32) also allows rigid-body motions. If  $\dot{\mathbf{a}} \neq \mathbf{0}$  corresponds to a rigid-body motion, this would create no strains in the body, i.e.  $\dot{\boldsymbol{\varepsilon}} = \mathbf{B}\dot{\mathbf{a}} = \mathbf{0}$ . It follows that the homogeneous equation system  $\mathbf{K}_t \dot{\mathbf{a}} = \mathbf{0}$  possesses a non-trivial solution  $\dot{\mathbf{a}} \neq \mathbf{0}$  and  $\mathbf{K}_t$  is therefore singular, i.e.  $\det \mathbf{K}_t = 0$ . This means, in fact, that we cannot solve the equation system (17.32).

To obtain a solution of our boundary value problem, we therefore have to introduce the pertinent boundary conditions that, among other things, ensure that rigid-body motions are prevented.

Let us write the standard iteration scheme (17.18) and (17.20) in a simplified form as

$$\mathbf{A} \delta \mathbf{a} = \mathbf{f} - \mathbf{f}_{int} \quad (17.37)$$

where

$$\mathbf{A} = \mathbf{A}^{i-1} ; \quad \delta \mathbf{a} = \mathbf{a}^i - \mathbf{a}^{i-1} \quad (17.38)$$

Moreover

$$\mathbf{f} = \int_S \mathbf{N}^T \mathbf{t}_{n+1} dS + \int_V \mathbf{N}^T \mathbf{b}_{n+1} dV ; \quad \mathbf{f}_{int} = \int_V \mathbf{B}^T \boldsymbol{\sigma}^{i-1} dV \quad (17.39)$$

where we notice that the stresses  $\boldsymbol{\sigma}^{i-1}$  are known, i.e. the internal forces  $\mathbf{f}_{int}$  are also known.

Since  $\mathbf{A}$  may always be interpreted as a stiffness matrix of the body, it follows in accordance with the discussion above that

$$\det \mathbf{A} = 0$$

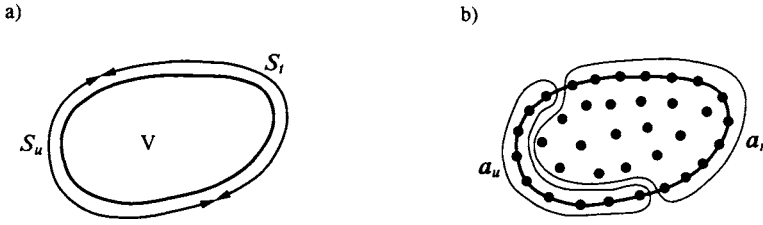
i.e., equation system (17.37) can first be solved after the boundary conditions have been introduced.

The boundary conditions are prescribed as

$\mathbf{u} = \text{is given along } S_u$ $\mathbf{t} = \text{is given along } S_t$
--

That is, the displacement vector  $\mathbf{u}$  is prescribed along the boundary surface  $S_u$  and the traction vector  $\mathbf{t}$  is prescribed along the boundary surface  $S_t$ . This is illustrated in Fig. 17.13a).

We now split the finite element nodes into two groups: one group along the boundary  $S_u$  where the displacements are prescribed and another group which contains the remaining nodes. Correspondingly, the nodal displacements  $\mathbf{a}$  can



**Figure 17.13:** a) Along  $S_u$  the displacements are given and along  $S_r$  the traction vector is known; b) split of nodal displacements into prescribed values  $a_u$  along  $S_u$  and the remaining, i.e. unknown, values  $a_r$  in the remaining part of the body.

then be partitioned into  $a_u$  and  $a_r$ , respectively, i.e.  $a_u$  are known whereas the remaining nodal displacements  $a_r$  are to be determined; this partitioning is illustrated in Fig. 17.13b). In accordance with this procedure, we may, in principle, partition (17.37) and with evident notation, we obtain

$$\begin{bmatrix} \mathbf{A}_{uu} & \mathbf{A}_{ur} \\ \mathbf{A}_{ru} & \mathbf{A}_{rr} \end{bmatrix} \begin{bmatrix} \delta \mathbf{a}_u \\ \delta \mathbf{a}_r \end{bmatrix} = \begin{bmatrix} \mathbf{f}_u - (\mathbf{f}_{int})_u \\ \mathbf{f}_r - (\mathbf{f}_{int})_r \end{bmatrix} \quad (17.40)$$

Since all the internal forces  $\mathbf{f}_{int}$  are known, both  $(\mathbf{f}_{int})_u$  and  $(\mathbf{f}_{int})_r$  are known forces. However, of the external forces  $\mathbf{f}$ , only the part  $\mathbf{f}_r$  is, in fact, known whereas the part  $\mathbf{f}_u$  is unknown, since it corresponds to the reactions where the displacements are prescribed.

From the second row of (17.40), we find

$$\mathbf{A}_{rr} \delta \mathbf{a}_r = \mathbf{f}_r - (\mathbf{f}_{int})_r - \mathbf{A}_{ru} \delta \mathbf{a}_u$$

where the right-hand side is known. Written explicitly, using (17.38) and (17.39), we have

$$\mathbf{A}_{rr}^{i-1} (\mathbf{a}_r^i - \mathbf{a}_r^{i-1}) = (\mathbf{f}_{n+1})_r - (\mathbf{f}_{int}^{i-1})_r - \mathbf{p}_r^{i-1} \quad (17.41)$$

where  $\mathbf{p}_r^{i-1}$  represents the quantity

$$\mathbf{p}_r^{i-1} = \mathbf{A}_{ru}^{i-1} (\mathbf{a}_u^i - \mathbf{a}_u^{i-1})$$

The known displacement changes along the boundary  $S_u$  are imposed in the very first integration, i.e.

$$\mathbf{p}_r^0 = (\mathbf{A}_{ru})_n [(\mathbf{a}_u)_{n+1} - (\mathbf{a}_u)_n]$$

In the next iterations, no change occurs of the displacements  $a_u$  along  $S_u$ , i.e.

$$\mathbf{p}_r^{i-1} = \mathbf{0} \quad \text{when} \quad i = 2, 3, \dots \quad (17.42)$$

In reality, it is on the equation system (17.41) we perform the equilibrium equations and all the previous as well as the following sections should be interpreted in this manner.

The iteration scheme will always take at least one iteration and it is only if the out-of-balance forces after this first iteration are too large that the iterations will continue. Therefore, we obtain with (17.41) and (17.42) that the part of the out-of-balance forces on which we measure equilibrium is  $\psi_r^{i-1} = (f_{int}^{i-1})_r - (f_{n+1})_r$ . Consequently

$$\boxed{\text{Fulfillment of equilibrium is checked for the out-of-balance forces } \psi_r^{i-1} = (f_{int}^{i-1})_r - (f_{n+1})_r} \quad (17.43)$$

The iterations are stopped when equilibrium has been fulfilled. In that case  $\delta a_r = a_r^i - a_r^{i-1} = 0$  and as also  $\delta a_u = 0$  (after the first iteration), (17.40) shows that the reaction forces along  $S_u$ , i.e.  $(f_{n+1})_u$  are given by

$$(f_{n+1})_u = (f_{int}^{i-1})_u \quad (17.44)$$

Since the stresses  $\sigma^{i-1}$  are known, the internal force  $f_{int}^{i-1}$  and thereby also the part  $(f_{int}^{i-1})_u$  are known. After equilibrium has been obtained, the unknown reactions  $(f_{n+1})_u$  along  $S_u$  are therefore determined by (17.44).

On the other hand, if we always determine the reaction forces by (17.44) then equilibrium along the boundary  $S_u$  is always fulfilled, i.e.  $\psi_u^{i-1} = (f_{int}^{i-1})_u - (f_{n+1})_u = 0$ , and as  $(\psi^{i-1})^T = [(\psi_r^{i-1})^T, (\psi_u^{i-1})^T]$ , we may instead of (17.43) check the equilibrium conditions on the total out-of-balance force  $\psi^{i-1}$ .

## 17.7 Convergence criteria

Let us return to the general iteration scheme (17.18)-(17.24) which starts with the iteration counter  $i = 1$  and then continues for increasing  $i$ -values. Hopefully, the iteration scheme converges in the sense that the out-of-balance forces  $\psi$  approach zero, which implies that the new solution  $a^i$  differs only insignificantly from the previous solution  $a^{i-1}$ . In practice however, the out-of-balance forces  $\psi$  will never be exactly zero, so we have to specify some threshold for  $\psi$  that terminates the iterations. Such a threshold value is called a *convergence criterion*.

Convergence criteria may be formulated in a number of different manners and for a more detailed discussion, we refer to Bathe and Cimento (1980), Bathe (1996), Bergan and Clough (1972) and Crisfield (1991). Instead of a threshold value for  $\psi$  we may measure the improvement of our solution in each iteration in terms of  $a^i - a^{i-1}$  and apply a convergence criterion for this quantity.

To enforce a threshold value on the out-of-balance forces  $\psi^{i-1}$ , we have to measure  $\psi^{i-1}$  against something and for this purpose it seems logical to choose

the external forces  $\mathbf{f}_{n+1}$ . The quantities  $\boldsymbol{\psi}^{i-1}$  and  $\mathbf{f}_{n+1}$  contain many components and in order to compare these quantities, we may use their lengths measured by the scalar product. Therefore, an often used force convergence criterion is given by

$$[(\boldsymbol{\psi}^{i-1})^T \boldsymbol{\psi}^{i-1}]^{1/2} \leq \varepsilon_{F_1} (\mathbf{f}_{n+1}^T \mathbf{f}_{n+1})^{1/2} \quad (17.45)$$

where  $\varepsilon_{F_1}$  expresses the threshold. Typical values for this threshold are  $\varepsilon_{F_1} = 10^{-3} - 10^{-2}$ , cf. Zienkiewicz and Taylor (1991) and Crisfield (1991), and the smaller the value  $\varepsilon_{F_1}$ , the less is the acceptable error in the out-of-balance forces. If  $\varepsilon_{F_1}$  is chosen too small, many (costly) iterations are performed without improving the solution significantly and if  $\varepsilon_{F_1}$  is chosen too large, the solution becomes inaccurate.

It is a general experience in the literature that it is difficult to recommend a specific tolerance that is always of relevance; one tolerance may work well in some case and may be inferior in other situations. For a given problem, one often starts with a rather crude tolerance and then performs calculations with narrower tolerances in an effort to judge whether a converged solution has been obtained.

In (17.45), the out-of-balance forces are compared with the total external force. When the load steps are small, the out-of-balance forces are often compared with the load step itself, i.e.

$$[(\boldsymbol{\psi}^{i-1})^T \boldsymbol{\psi}^{i-1}]^{1/2} \leq \varepsilon_{F_2} [\Delta \mathbf{f}^T \Delta \mathbf{f}]^{1/2}$$

where  $\Delta \mathbf{f} = \mathbf{f}_{n+1} - \mathbf{f}_n$ ; Bathe and Cimento (1980) suggested  $\varepsilon_{F_2} \approx 0.1$ .

For plastic structures with a small plastic modulus  $H$ , force criteria may be misleading; the out-of-balance force may be small, but the structure is far from its correct displacement pattern. Therefore, in addition to the force convergence criteria, use is often made of some displacement convergence criteria, for instance

$$[(\delta \mathbf{u}^i)^T \delta \mathbf{u}^i]^{1/2} \leq \varepsilon_D [\mathbf{a}_n^T \mathbf{a}_n]^{1/2}$$

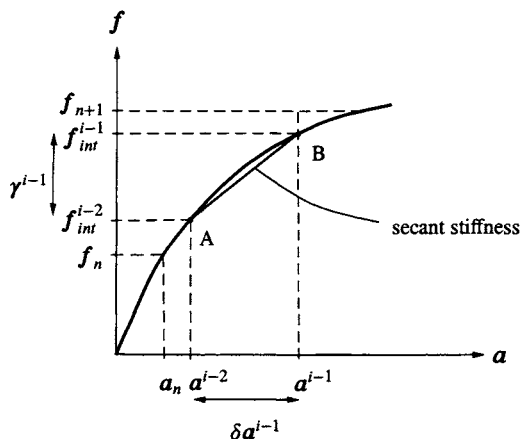
where  $\delta \mathbf{u}^i = \mathbf{a}^i - \mathbf{a}^{i-1}$  and  $\varepsilon_D \approx 10^{-3}$ , cf. Bathe and Cimento (1980).

The convergence criteria above possess the drawback that quantities having one dimension may be added to quantities having another dimension; the components of  $\mathbf{a}$  may involve both displacements and rotations and, likewise, the components of  $\boldsymbol{\psi}$  may involve forces as well as moments. A convergence criterion in terms of an energy criterion avoids this problem since force components are here multiplied by displacements whereas moments are multiplied by rotations, i.e. it involves so-called *conjugated quantities*. Therefore an energy convergence criterion in the form

$$[(\boldsymbol{\psi}^{i-1})^T \mathbf{a}_n]^{1/2} \leq \varepsilon_E [\mathbf{f}_n^T \mathbf{a}_n]^{1/2}$$

is often adopted; the tolerance  $\varepsilon_E = 10^{-3} - 10^{-2}$  is often used, cf. Zienkiewicz and Taylor (1991).





**Figure 17.14:** Current state given by point B and previous state given by point A; secant stiffness between these points.

## 17.8 Quasi-Newton methods

The Newton-Raphson scheme possesses a number of advantages among which its fast convergence is the major reason for its popularity, cf. (17.34). However, it was mentioned that every Newton-Raphson iteration requires the formulation of a new tangent stiffness matrix  $K_t$  and that, in principle, the inverse matrix  $K_t^{-1}$  also needs to be established; this means that every Newton-Raphson iteration is costly. This drawback can be obviated by the initial stiffness method where the elastic stiffness matrix is used in every iteration or in the modified Newton-Raphson method where the stiffness matrix is occasionally updated. However, these approaches have a slower convergence rate than the full Newton-Raphson scheme. It is of considerable interest that there exists an entirely different approach by which we can maintain a fast convergence without performing a costly matrix inversion in every iteration.

In the previous iteration schemes, we only took advantage of the information of the current state defined by  $a^{i-1}$  and  $\psi^{i-1}$  and with this information, we tried to establish a new estimate  $a^i$  and thereby  $\psi^i$ . Essentially, this means that all previous information of the response is ignored. To illustrate this viewpoint, consider Fig. 17.14 and suppose that we increase the external loading from  $f_n$  to  $f_{n+1}$ . Equilibrium iterations are then performed with, say, a Newton-Raphson procedure. At the present stage, we have then obtained the nodal displacements  $a^{i-1}$  at point B whereas the previous nodal displacements are given by  $a^{i-2}$  at point A.

In a Newton-Raphson procedure, we then determine the tangential stiffness matrix  $K_t$  at point B to obtain the new nodal displacement estimate  $a^i$ . Our objective is now to establish a procedure that as closely as possible retains the

fast convergence of the Newton-Raphson approach. Instead of the tangential stiffness at point B, it is then tempting to use the secant stiffness between point A and point B. This secant stiffness is a close approximation to the tangential stiffness at point B. To define this secant stiffness matrix, we observe that the internal forces at point B are given by  $f_{int}^{i-1} = \int_V \mathbf{B}^T \boldsymbol{\sigma}^{i-1} dV$  whereas the internal forces at point A are given by  $f_{int}^{i-2} = \int_V \mathbf{B}^T \boldsymbol{\sigma}^{i-2} dV$ . These internal forces are illustrated in Fig. 17.14. Define the quantities  $\delta \mathbf{a}^{i-1}$  and  $\boldsymbol{\gamma}^{i-1}$  by

$$\delta \mathbf{a}^{i-1} = \mathbf{a}^{i-1} - \mathbf{a}^{i-2}; \quad \boldsymbol{\gamma}^{i-1} = f_{int}^{i-1} - f_{int}^{i-2} \quad (17.46)$$

These quantities are illustrated in Fig. 17.14. The secant stiffness  $\mathbf{K}_s^{i-1}$  between point A and B is then defined by

$$\mathbf{K}_s^{i-1} \delta \mathbf{a}^{i-1} = \boldsymbol{\gamma}^{i-1}$$

which implies

$$\boxed{\delta \mathbf{a}^{i-1} = \mathbf{H}^{i-1} \boldsymbol{\gamma}^{i-1} \quad \text{Quasi-Newton relation}} \quad (17.47)$$

where  $\mathbf{H}^{i-1}$  is the inverse of  $\mathbf{K}_s^{i-1}$ . This relation is called the *quasi-Newton relation* since it defines a matrix  $\mathbf{H}^{i-1}$  that is almost equal to the inverse of the tangential stiffness at point B that arises in the Newton-Raphson scheme. It appears that if in the general iteration scheme (17.18) we make the choice

$$(\mathbf{A}^{i-1})^{-1} = \mathbf{H}^{i-1} \quad (17.48)$$

then we will obtain a convergence that is almost as fast as the Newton-Raphson scheme.

Before we proceed further, assume that the total number of degrees of freedom in the finite element discretization is  $N$ , then the matrix  $\mathbf{H}^{i-1}$  has the dimension  $N \times N$ . However, the quasi-Newton equation (17.47) - i.e. the secant relation - only comprises  $N$  equations and this implies that  $\mathbf{H}^{i-1}$  is not defined uniquely by (17.47) unless in the trivial case where  $N = 1$ . To establish  $\mathbf{H}^{i-1}$  so that it satisfies (17.47) therefore leaves us with considerable freedom.

That apart, our goal was to establish a method by which costly matrix inversions are avoided. Therefore, suppose that in the previous iteration we have obtained the matrix  $\mathbf{H}^{i-2}$  we then want to be able to establish  $\mathbf{H}^{i-1}$  by means of the simple updating scheme

$$\boxed{\mathbf{H}^{i-1} = \mathbf{H}^{i-2} + \mathbf{H}_{corr}^{i-1}} \quad (17.49)$$

where the *correction matrix*  $\mathbf{H}_{corr}^{i-1}$  should be easy to determine. If this can be achieved, we have fulfilled our objective of devising a method with almost as good convergence properties as the Newton-Raphson method without having to perform costly inversions of some stiffness matrix.

Methods that satisfy (17.47)-(17.49) are called *quasi-Newton methods* - occasionally also called *variable matrix methods* - and they were originally derived within nonlinear optimization theory. The literature on quasi-Newton methods is very extensive and reviews are given by Dennis and Moré (1977), Fletcher (1980) and Luenberger (1984).

### 17.8.1 Rank one correction

To provide a simple illustration of a quasi-Newton method, suppose that we take the updating scheme (17.49) in the form

$$\mathbf{H}^{i-1} = \mathbf{H}^{i-2} + \mathbf{v}^{i-1}(\mathbf{v}^{i-1})^T \quad (17.50)$$

where  $\mathbf{v}^{i-1}$  is a column matrix that is to be determined so that the quasi-Newton equation (17.47) is fulfilled. Insertion of (17.50) in (17.47) gives

$$\delta \mathbf{a}^{i-1} = \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1} + \mathbf{v}^{i-1}((\mathbf{v}^{i-1})^T \boldsymbol{\gamma}^{i-1})$$

which provides

$$\mathbf{v}^{i-1} = \frac{1}{(\mathbf{v}^{i-1})^T \boldsymbol{\gamma}^{i-1}} (\delta \mathbf{a}^{i-1} - \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1}) \quad (17.51)$$

and thereby

$$((\mathbf{v}^{i-1})^T \boldsymbol{\gamma}^{i-1})^2 = (\delta \mathbf{a}^{i-1} - \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1})^T \boldsymbol{\gamma}^{i-1}$$

Insertion of (17.51) in (17.50) and use of the expression above give the result

$$\mathbf{H}^{i-1} = \mathbf{H}^{i-2} + \frac{(\delta \mathbf{a}^{i-1} - \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1})(\delta \mathbf{a}^{i-1} - \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1})^T}{(\delta \mathbf{a}^{i-1} - \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1})^T \boldsymbol{\gamma}^{i-1}} \quad (17.52)$$

Since  $\delta \mathbf{a}^{i-1}$  and  $\boldsymbol{\gamma}^{i-1}$  are known, it appears that if  $\mathbf{H}^{i-2}$  is also known then the new value of  $\mathbf{H}^{i-1}$ , i.e. the new value of the inverse of the secant stiffness matrix, is obtained by performing the simple operations in (17.52) and no costly inversion technique is necessary.

To further illustrate this approach, let us recall some results from matrix algebra, cf. for instance Strang (1980). For any square matrix  $\mathbf{P}$ , the eigenvalue problem is given by  $\mathbf{P}\mathbf{c} = \lambda\mathbf{c}$ , where  $\mathbf{c}$  is the eigenvector and  $\lambda$  is the corresponding eigenvalue. The eigenvalue problem can also be written as the homogeneous equation system  $(\mathbf{P} - \lambda\mathbf{I})\mathbf{c} = \mathbf{0}$  and a non-trivial  $\mathbf{c}$ -solution requires that  $\det(\mathbf{P} - \lambda\mathbf{I}) = 0$ ; this condition is the *characteristic equation*. If the matrix  $\mathbf{P}$  possesses an eigenvalue that is zero, then we must have  $\det \mathbf{P} = 0$ . Let the matrix  $\mathbf{P}$  be of dimension  $N \times N$ , then the characteristic equation  $\det(\mathbf{P} - \lambda\mathbf{I}) = 0$  becomes a polynomial of order  $N$  and we therefore have  $N$  eigenvalues. The number of non-zero eigenvalues of the matrix  $\mathbf{P}$  is called the *rank of the matrix* and it is denoted by  $r$ . If  $\det(\mathbf{P} - \lambda\mathbf{I}) = 0$  possesses no zero eigenvalues, the

rank is  $r = N$ . If it possesses one zero eigenvalue, the rank is  $r = N - 1$  and if it possesses two zero eigenvalues, the rank is  $r = N - 2$  and it is then said that the matrix  $\mathbf{P}$  has the eigenvalue zero with a *multiplicity* of two. Therefore, if the homogeneous equation system  $\mathbf{P}\mathbf{c} = \mathbf{0}$  has  $q$  non-trivial  $\mathbf{c}$  - solutions that are linearly independent then the rank of  $\mathbf{P}$  is  $r = N - q$  and the matrix  $\mathbf{P}$  possesses the zero eigenvalue with a multiplicity of  $q$ .

In the present case, consider the correction matrix  $\mathbf{H}_{corr}^{i-1}$  which according to (17.49) and (17.50) is given by  $\mathbf{H}_{corr}^{i-1} = \mathbf{v}^{i-1}(\mathbf{v}^{i-1})^T$  and the dimension of  $\mathbf{H}_{corr}^{i-1}$  is  $N \times N$ . Investigate the possibility of a non-trivial solution of the homogeneous equation system  $\mathbf{H}_{corr}^{i-1}\mathbf{c} = \mathbf{0}$ , i.e.  $\mathbf{v}^{i-1}((\mathbf{v}^{i-1})^T\mathbf{c}) = \mathbf{0}$ . It appears that there are  $q = N - 1$  linear independent  $\mathbf{c}$ -vectors that fulfill  $(\mathbf{v}^{i-1})^T\mathbf{c} = 0$ , i.e.  $\mathbf{H}_{corr}^{i-1}$  has a zero eigenvalue with a multiplicity of  $q = N - 1$ . The rank of  $\mathbf{H}_{corr}^{i-1}$  is therefore  $r = N - q = N - (N - 1) = 1$ . Therefore, the quasi-Newton method defined by (17.50) is called a *rank one correction*.

The rank one correction quasi-Newton method defined by (17.50) was suggested by Broyden (1967) and later independently by others. It appears that if  $\mathbf{H}^{i-2}$  is symmetric then this method has the neat property that the new update  $\mathbf{H}^{i-1}$  is also symmetric. However, if the denominator in (17.52) is also small this may lead to numerical difficulties. Therefore, we shall not dwell on this approach any further and here merely take it as a simple illustration of a quasi-Newton method.

### 17.8.2 Rank two correction. BFGS-method

To establish a quasi-Newton method that is of relevance for our purpose, consider the following form

$$\mathbf{H}^{i-1} = (\mathbf{I} + \mathbf{v}\mathbf{w}^T)^T \mathbf{H}^{i-2} (\mathbf{I} + \mathbf{v}\mathbf{w}^T) \quad (17.53)$$

where the matrix  $\mathbf{v}\mathbf{w}^T$  must be chosen so that (17.53) fulfills the quasi-Newton equation (17.47). Carrying out the multiplications and assuming  $\mathbf{H}^{i-2}$  to be symmetric we obtain

$$\mathbf{H}^{i-1} = \mathbf{H}^{i-2} + (\mathbf{v}^T \mathbf{H}^{i-2} \mathbf{v}) \mathbf{w}\mathbf{w}^T + \mathbf{w}(\mathbf{H}^{i-2} \mathbf{v})^T + \mathbf{H}^{i-2} \mathbf{v}\mathbf{w}^T \quad (17.54)$$

which is evidently of the format given by (17.49). In the quasi-Newton approach defined by (17.50), the correction matrix  $\mathbf{H}_{corr}^{i-1}$  was found to be of rank one. Let us now determine the rank of the correction matrix in the approach defined by (17.54). Therefore, we shall investigate non-trivial  $\mathbf{c}$ -solutions of the homogeneous equation system  $\mathbf{H}_{corr}^{i-1}\mathbf{c} = \mathbf{0}$ . With (17.54) and (17.49) and defining the vector  $\mathbf{b}$  by  $\mathbf{b} = \mathbf{H}^{i-2}\mathbf{v}$ , the following equation system is therefore considered

$$\mathbf{H}_{corr}^{i-1}\mathbf{c} = (\mathbf{v}^T \mathbf{b}) \mathbf{w}(\mathbf{w}^T \mathbf{c}) + \mathbf{w}(\mathbf{b}^T \mathbf{c}) + \mathbf{b}(\mathbf{w}^T \mathbf{c}) = \mathbf{0}$$

This equation system is fulfilled if both  $\mathbf{w}^T \mathbf{c} = 0$  and  $\mathbf{b}^T \mathbf{c} = 0$ . With  $\mathbf{H}_{corr}^{i-1}$  having the dimension  $N \times N$ , these conditions can be fulfilled for  $N - 2$  linear

independent  $c$ -vectors. Referring to our previous discussion,  $\mathbf{H}_{corr}^{i-1}$  has therefore the eigenvalue zero with a multiplicity of  $q = N - 2$  and the rank  $r$  of  $\mathbf{H}_{corr}^{i-1}$  is then  $r = N - q = N - (N - 2) = 2$ . Consequently, the quasi-Newton approach defined by (17.54) is a *rank two correction*.

The next topic is to identify the vectors  $\mathbf{v}$  and  $\mathbf{w}$  in (17.54). Since the vector  $\delta\mathbf{a}^{i-1}$  is one of the vectors we know, it seems natural to choose  $\mathbf{w}$  to be proportional to  $\delta\mathbf{a}^{i-1}$ . According to (17.53),  $\mathbf{v}$  and  $\mathbf{w}$  only occur in the combination  $\mathbf{v}\mathbf{w}^T$  or  $\mathbf{w}\mathbf{v}^T$  and the length of one of these vectors can therefore be chosen arbitrarily; this will only affect the length of the other vector. Let us therefore choose  $\mathbf{w}$  according to

$$\mathbf{w} = \frac{1}{b} \delta\mathbf{a}^{i-1} \quad \text{where} \quad b = (\boldsymbol{\gamma}^{i-1})^T \delta\mathbf{a}^{i-1} \quad (17.55)$$

Since both  $\delta\mathbf{a}^{i-1}$  and  $\boldsymbol{\gamma}^{i-1}$  are known, also  $\mathbf{w}$  is known.

With (17.55), (17.54) becomes

$$\begin{aligned} \mathbf{H}^{i-1} &= \mathbf{H}^{i-2} + (\mathbf{v}^T \mathbf{H}^{i-2} \mathbf{v}) \frac{1}{b^2} \delta\mathbf{a}^{i-1} (\delta\mathbf{a}^{i-1})^T + \frac{1}{b} \delta\mathbf{a}^{i-1} (\mathbf{H}^{i-2} \mathbf{v})^T \\ &\quad + \frac{1}{b} \mathbf{H}^{i-2} \mathbf{v} (\delta\mathbf{a}^{i-1})^T \end{aligned} \quad (17.56)$$

Insertion into the quasi-Newton equation (17.47) results in

$$\begin{aligned} \delta\mathbf{a}^{i-1} &= \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1} + (\mathbf{v}^T \mathbf{H}^{i-2} \mathbf{v}) \frac{1}{b} \delta\mathbf{a}^{i-1} \\ &\quad + \frac{1}{b} \delta\mathbf{a}^{i-1} (\mathbf{v}^T \mathbf{H}^{i-2} \boldsymbol{\gamma}^{i-1}) + \mathbf{H}^{i-2} \mathbf{v} \end{aligned}$$

where it was assumed that  $\mathbf{H}^{i-2}$  is symmetric. Determination of  $\mathbf{v}$  from the expression gives

$$\mathbf{v} = m(\mathbf{H}^{i-2})^{-1} \delta\mathbf{a}^{i-1} - \boldsymbol{\gamma}^{i-1} \quad (17.57)$$

where

$$m = 1 - \frac{1}{b} \mathbf{v}^T \mathbf{H}^{i-2} (\mathbf{v} + \boldsymbol{\gamma}^{i-1}) \quad (17.58)$$

Expression (17.57) leads to

$$\mathbf{v}^T \mathbf{H}^{i-2} (\mathbf{v} + \boldsymbol{\gamma}^{i-1}) = m^2 (\delta\mathbf{a}^{i-1})^T (\mathbf{H}^{i-2})^{-1} \delta\mathbf{a}^{i-1} - m b \quad (17.59)$$

where advantage was taken of (17.55). Insertion of (17.59) in (17.58) results in

$$m = \pm \left[ \frac{b}{(\delta\mathbf{a}^{i-1})^T (\mathbf{H}^{i-2})^{-1} \delta\mathbf{a}^{i-1}} \right]^{1/2} \quad (17.60)$$

The vector  $\mathbf{v}$  given by (17.57) and (17.60) has now been expressed in terms of known quantities. Moreover, this expression for  $\mathbf{v}$  ensures that  $\mathbf{H}^{i-1}$  given

by (17.56) fulfills the quasi-Newton equation (17.47). It is easy to be convinced that irrespective of the sign we choose for  $m$  in (17.60), (17.56) provides the same updated matrix  $H^{i-1}$ . In (17.60), we shall for convenience choose the positive sign.

In view of this comment, the results (17.53), (17.55), (17.57) and (17.60) can be summarized as

<i>BFGS-method</i>	
$H^{i-1} = (I + \mathbf{v}\mathbf{w}^T)^T H^{i-2} (I + \mathbf{v}\mathbf{w}^T)$	
where	
$\mathbf{v} = \left[ \frac{(\gamma^{i-1})^T \delta \mathbf{a}^{i-1}}{(\delta \mathbf{a}^{i-1})^T (H^{i-2})^{-1} \delta \mathbf{a}^{i-1}} \right]^{1/2}$	(17.61)
and	
$\mathbf{w} = \frac{\delta \mathbf{a}^{i-1}}{(\gamma^{i-1})^T \delta \mathbf{a}^{i-1}}$	

This comprises the so-called *BFGS-method*, suggested by Broyden (1970), Fletcher (1970), Goldfarb (1970) and Shannon (1970). The BFGS-method can be expressed in various identical forms and the product form adopted in (17.61) was proposed by Matthies and Strang (1979). Indeed, Matthies and Strang (1979) were the first to apply the BFGS-method in a finite element context.

In the BFGS-scheme (17.61), we apparently need to establish the inverse matrix  $(H^{i-2})^{-1}$ . This inverse matrix always occurs in the form  $(H^{i-2})^{-1} \delta \mathbf{a}^{i-1}$  and we will now show that this term is already known. From (17.46a) and the general iteration scheme (17.18), we have

$$\mathbf{A}^{i-2} \delta \mathbf{a}^{i-1} = \mathbf{A}^{i-2} (\mathbf{a}^{i-1} - \mathbf{a}^{i-2}) = -\boldsymbol{\psi}^{i-2}$$

and since (17.48) shows that  $\mathbf{A}^{i-2} = (\mathbf{H}^{i-2})^{-1}$ , we obtain

$$(\mathbf{H}^{i-2})^{-1} \delta \mathbf{a}^{i-1} = -\boldsymbol{\psi}^{i-2}$$

The out-of-balance force  $\boldsymbol{\psi}^{i-2}$  is already known and this facilitates the use of the BFGS-scheme (17.61) considerably.

The product form for the updated matrix  $H^{i-1}$  given by (17.61) is of advantage for several reasons. Defining the matrix  $\mathbf{M}$  by

$$\mathbf{M} = I + \mathbf{v}\mathbf{w}^T \quad (17.62)$$

we obtain

$$H^{i-1} = \mathbf{M}^T H^{i-2} \mathbf{M} \quad (17.63)$$

The stiffness matrix in the finite element formulation is often symmetric and so is therefore its inverse. In the previous derivations, we assumed  $H^{i-2}$  to be

symmetric and (17.63) then shows the interesting result that the updated matrix  $\mathbf{H}^{i-1}$  is also symmetric. Let us next assume that  $\mathbf{H}^{i-2}$  is positive definite. To investigate whether the updated matrix  $\mathbf{H}^{i-1}$  is also positive definite, we choose an arbitrary vector  $\mathbf{c}$  and calculate the quadratic form of  $\mathbf{H}^{i-1}$ . Denoting the value of the quadratic form by  $a$ , we obtain with (17.63)

$$a = \mathbf{c}^T \mathbf{H}^{i-1} \mathbf{c} = \mathbf{y}^T \mathbf{H}^{i-2} \mathbf{y} \quad (17.64)$$

where  $\mathbf{y}$  is defined by

$$\mathbf{M} \mathbf{c} = \mathbf{y}$$

It appears that if  $\mathbf{c} \neq \mathbf{0}$  implies  $\mathbf{y} \neq \mathbf{0}$  then, since  $\mathbf{H}^{i-2}$  was assumed to be positive definite,  $a > 0$  and  $\mathbf{H}^{i-1}$  is therefore also positive definite.

The only situation where  $\mathbf{c} \neq \mathbf{0}$  implies  $\mathbf{y} = \mathbf{0}$  is if  $\mathbf{M} \mathbf{c} = \mathbf{0}$  possesses a non-trivial solution. In turn this requires that the matrix  $\mathbf{M}$  has an eigenvalue that is zero. Let us therefore determine the eigenvalues  $\lambda$  of the matrix  $\mathbf{M}$ , i.e.  $\mathbf{M} \mathbf{c} = \lambda \mathbf{c}$ . With (17.62), we then obtain

$$(1 - \lambda) \mathbf{c} + \mathbf{v}(\mathbf{w}^T \mathbf{c}) = \mathbf{0} \quad (17.65)$$

Let us assume that  $\mathbf{w}^T \mathbf{c} = 0$ ; with the dimension of  $\mathbf{M}$  being  $N \times N$ , this condition is fulfilled by  $N - 1$  linearly independent  $\mathbf{c}$ -vectors and for each of these  $\mathbf{c}$ -vectors, (17.65) is fulfilled for  $\lambda = 1$ . We conclude

$$\lambda = 1 \quad \text{with multiplicity } N - 1$$

Assume next that  $\mathbf{w}^T \mathbf{c} \neq 0$ . Expression (17.65) then shows that the eigenvector  $\mathbf{c}$  must be proportional to  $\mathbf{v}$ , i.e.  $\mathbf{c} = \alpha \mathbf{v}$  where  $\alpha$  is a constant. Use of  $\mathbf{c} = \alpha \mathbf{v}$  in (17.65) gives  $[1 - \lambda + (\mathbf{w}^T \mathbf{v})] \mathbf{v} = \mathbf{0}$ , i.e. we obtain the following remaining eigenvalue

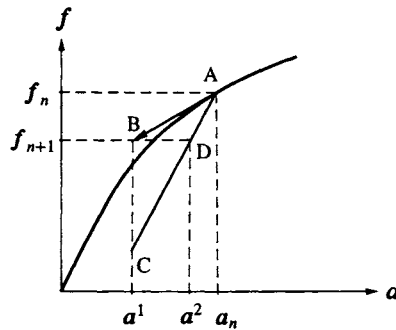
$$\lambda = 1 + (\mathbf{w}^T \mathbf{v})$$

which with  $\mathbf{v}$  and  $\mathbf{w}$  given by (17.61) becomes

$$\lambda = \left[ \frac{(\delta \mathbf{a}^{i-1})^T (\mathbf{H}^{i-2})^{-1} \delta \mathbf{a}^{i-1}}{(\mathbf{y}^{i-1})^T \delta \mathbf{a}^{i-1}} \right]^{1/2} \quad (17.66)$$

Since  $\mathbf{H}^{i-2}$  and thereby  $(\mathbf{H}^{i-2})^{-1}$  was assumed to be positive definite, the eigenvalue given by (17.66) can never be zero. In conclusion, we have proved that the matrix  $\mathbf{M}$  has no zero eigenvalues and the discussion relating to (17.64) then shows that the updated matrix  $\mathbf{H}^{i-1}$  is positive definite, if  $\mathbf{H}^{i-2}$  is so. In conclusion

*If  $\mathbf{H}^{i-2}$  is symmetric and positive definite then the BFGS-method implies that the updated matrix  $\mathbf{H}^{i-1}$  is symmetric and positive definite*



**Figure 17.15:** BFGS-method during unloading of elasto-plastic body. Correct solution is obtained after two iterations.

This interesting property implies, for instance, that if the very first  $\mathbf{H}^{i-2}$ -matrix is chosen as the inverse of the elastic stiffness matrix, which certainly is symmetric and positive definite, then all the following updates  $\mathbf{H}^{i-1}$  generated by (17.61) are symmetric and positive definite.

We will now show that the BFGS-method for unloading of a elasto-plastic body provides the correct elastic unloading after two iterations; this important property is a result of the secant relation expressed by the quasi-Newton equation (17.47). In the first iteration and referring to Fig. 17.15, unloading takes place along some secant direction AB giving the first displacement estimate  $a^1$ . The internal forces corresponding to  $a^1$  are given by point C. In the next iteration, the secant stiffness is determined by the last two sets of displacements and the last two sets of internal forces. Since we can assume that equilibrium is fulfilled at point A, this new secant stiffness is given by CA equal to the elastic stiffness. With this stiffness together with the residual force vector given by CB we reach, in the second iteration, exactly the correct response given by point D.

Efficient implementation of the BFGS-method evidently hinges on keeping the number of computer operations as small as possible and this requires a number of considerations. As we are here only concerned with the concepts behind the BFGS-method, the reader is referred to Bathe and Cimento (1980), Crisfield (1991), Fletcher (1980), Luenberger (1984) and Matthies and Strang (1979) for information on implementation techniques. Here we merely observe that the BFGS-method provides an attractive approach, which gives a convergence rate that is almost as fast as the Newton-Raphson without performing costly matrix inversions. The BFGS-method therefore provides an interesting compromise between the Newton-Raphson scheme and the modified Newton-Raphson scheme. The advantages within a finite element context are well documented, see Bathe (1996), Bathe and Cimento (1980), Crisfield (1991) and Matthies and Strang (1979). The use of a nonsymmetric quasi-Newton approach in finite element calculations of fluid mechanics has been investigated by Engelman *et al.* (1981).



## 17.9 Line search

Irrespective of the choice of iteration matrix  $\mathbf{A}$ , it turns out to be possible to reduce the number of equilibrium iterations significantly by introducing the concept of *line search*. Like many other issues relating to solution of nonlinear equations, this concept has its origin in optimization theory, see for instance Luenberger (1984).

Let us define the quantity  $\mathbf{s}^{i-1}$  by

$$\boxed{\mathbf{s}^{i-1} = -(\mathbf{A}^{i-1})^{-1} \boldsymbol{\psi}^{i-1} \quad \text{search direction}} \quad (17.67)$$

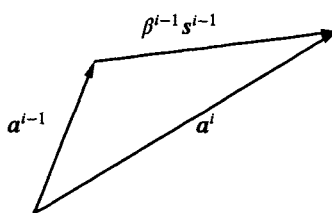
Then the general iteration scheme (17.18) can be written as

$$\mathbf{a}^i = \mathbf{a}^{i-1} + \mathbf{s}^{i-1} \quad (17.68)$$

It appears that the new solution  $\mathbf{a}^i$  is obtained by correcting the old solution  $\mathbf{a}^{i-1}$  by the quantity  $\mathbf{s}^{i-1}$ . This quantity is called the *search direction* since it is the direction in which we search for the new solution. With this interpretation, it is tempting to accept the search direction  $\mathbf{s}^{i-1}$ , but leave open how far we shall go in this direction. Therefore, instead of (17.68), we now adopt the scheme

$$\boxed{\mathbf{a}^i = \mathbf{a}^{i-1} + \beta^{i-1} \mathbf{s}^{i-1}} \quad (17.69)$$

where the parameter  $\beta^{i-1}$ , the *acceleration factor*, is to be determined. With this viewpoint, we accept the search direction  $\mathbf{s}^{i-1}$  and adjust the acceleration factor  $\beta^{i-1}$  in some way so that (17.69) gives us the best possible new estimate. The scheme is illustrated in Fig. 17.16



**Figure 17.16:** Scheme when line search is adopted

A comparison of (17.69) and (17.18) shows that instead of working with the iteration matrix  $\mathbf{A}^{i-1}$  we now work with the iteration matrix  $\frac{1}{\beta^{i-1}} \mathbf{A}^{i-1}$ . The new scheme is therefore fully acceptable as long as  $\beta^{i-1} \neq 0$  and it appears that line search can be applied in combination with any of the methods discussed previously.

However, the question remains of how to determine  $\beta^{i-1}$  - occasionally called the *step length* - so that the scheme (17.69) is optimal in some sense. The concept of *line search* means that we go from an old solution  $\mathbf{a}^{i-1}$  corresponding to  $\beta^{i-1} = 0$  in the search direction  $\mathbf{s}^{i-1}$  until we obtain the best possible new solution  $\mathbf{a}^i$ . Since the iteration scheme without line search corresponds to  $\beta^{i-1} = 1$ , it is evident that we must have  $\beta^{i-1} > 0$ . The best estimate for  $\mathbf{a}^i$  is the one which makes the out-of-balance forces  $\boldsymbol{\psi}(\mathbf{a}^i) = \mathbf{0}$ . In general, this condition cannot be fulfilled since it requires that the search direction be the correct one. This is just to say that we cannot expect to obtain that all components of  $\boldsymbol{\psi}(\mathbf{a}^i)$  become zero just by adjusting one parameter  $\beta^{i-1}$ . Since  $\beta^{i-1}$  is one parameter, we must establish one condition from which  $\beta^{i-1}$  can be determined. Such a condition can be obtained by requiring that the new out-of-balance force  $\boldsymbol{\psi}(\mathbf{a}^i)$  is orthogonal to a vector, i.e. the component of  $\boldsymbol{\psi}(\mathbf{a}^i)$  in the direction of this vector is required to be zero. The question then arises which vector we shall choose when enforcing this orthogonality requirement. Referring to Fig. 17.16, three vectors are possible candidates:  $\mathbf{a}^i$ ,  $\mathbf{a}^{i-1}$  and  $\mathbf{s}^{i-1}$ . It turns out that the search direction  $\mathbf{s}^{i-1}$  is the most natural choice, i.e. we adopt the orthogonality condition  $(\mathbf{s}^{i-1})^T \boldsymbol{\psi}(\mathbf{a}^i) = 0$ , which with (17.69) reads

$$(\mathbf{s}^{i-1})^T \boldsymbol{\psi}(\mathbf{a}^{i-1} + \beta^{i-1} \mathbf{s}^{i-1}) = 0 \quad \text{line search} \quad (17.70)$$

Since  $\mathbf{s}^{i-1}$  and  $\mathbf{a}^{i-1}$  are known quantities, this orthogonality condition provides one (nonlinear) relation by which  $\beta^{i-1}$  can be determined.

To motivate the use of the search direction  $\mathbf{s}^{i-1}$  in the orthogonality condition (17.70), it is instructive to consider linear elasticity. In that case  $\boldsymbol{\sigma} = \mathbf{D}\boldsymbol{\epsilon} = \mathbf{D}\mathbf{B}\mathbf{a}$  and the out-of-balance forces  $\boldsymbol{\psi}$  defined by (17.2) then becomes

$$\boldsymbol{\psi} = \mathbf{K}\mathbf{a} - \mathbf{f} \quad (17.71)$$

where  $\mathbf{K} = \int_V \mathbf{B}^T \mathbf{D} \mathbf{B} dV$  is the elastic stiffness matrix. The equilibrium condition  $\boldsymbol{\psi} = \mathbf{0}$  is then fulfilled by solving the linear equation system  $\mathbf{K}\mathbf{a} - \mathbf{f} = \mathbf{0}$ .

Further insight into the linear problem can be obtained by introducing the *potential energy*  $\Pi$  of the linear elastic body. If the body is discretized by means of finite elements, the potential energy  $\Pi$  is defined by

$$\Pi = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{f} \quad (17.72)$$

where the external force  $\mathbf{f}$  is viewed as constant. It appears that  $\Pi = \Pi(\mathbf{a})$  and therefore

$$\frac{\partial \Pi}{\partial \mathbf{a}} = \mathbf{K}\mathbf{a} - \mathbf{f}$$

A comparison with (17.71) shows that equilibrium is expressed as  $\partial \Pi / \partial \mathbf{a} = \mathbf{0}$ , i.e. at equilibrium the potential energy  $\Pi$  takes an extremum value.

Evidently, it is possible to solve the linear equation system  $\mathbf{K}\mathbf{a} - \mathbf{f} = \mathbf{0}$  directly, but it is also possible to solve this linear equation system in an iterative manner using the scheme (17.69). With  $\mathbf{a}^i$  given by (17.69), the potential energy defined by (17.72) becomes.

$$\Pi(\mathbf{a}^i) = \frac{1}{2}(\mathbf{a}^i)^T \mathbf{K} \mathbf{a}^i - (\mathbf{a}^i)^T \mathbf{f} \quad (17.73)$$

Since  $\mathbf{a}^{i-1}$  and  $\mathbf{s}^{i-1}$  are given quantities,  $\mathbf{a}^i$  given by (17.69) depends only on  $\beta^{i-1}$  and we therefore have  $\Pi = \Pi(\beta^{i-1})$ . From (17.73) it then follows that

$$\frac{d\Pi}{d\beta^{i-1}} = (\mathbf{s}^{i-1})^T (\mathbf{K} \mathbf{a}^i - \mathbf{f}) \quad (17.74)$$

and

$$\frac{d^2\Pi}{d(\beta^{i-1})^2} = (\mathbf{s}^{i-1})^T \mathbf{K} \mathbf{s}^{i-1} > 0 \quad (17.75)$$

where it was used that  $\mathbf{K}$  is positive definite. In accordance with (17.71), we have for linear elasticity that

$$\boldsymbol{\psi}^i = \mathbf{K} \mathbf{a}^i - \mathbf{f} \quad (17.76)$$

We observed above that the potential energy takes an extremum value at equilibrium. As we are performing equilibrium iterations, the condition of equilibrium has not been achieved, but with the quantities  $\mathbf{a}^{i-1}$  and  $\mathbf{s}^{i-1}$  given, the best possible new solution is obtained by requiring  $\Pi$  to be extremum, i.e.  $d\Pi/d\beta^{i-1} = 0$ . Therefore, (17.74) gives with (17.76)

$$(\mathbf{s}^{i-1})^T \boldsymbol{\psi}^i = 0 \quad (17.77)$$

According to (17.75),  $d^2\Pi/d(\beta^{i-1})^2 > 0$  and the extremum property of the potential energy  $\Pi$  is therefore that of a minimum. Similar to (17.76), we have

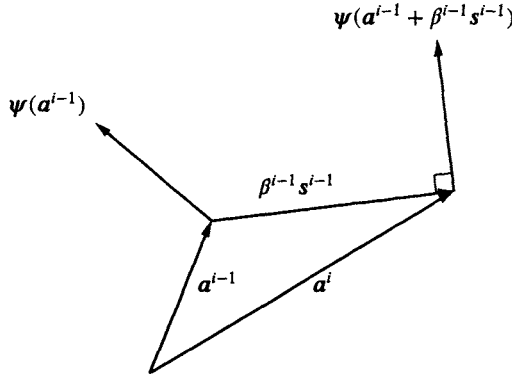
$$\boldsymbol{\psi}^{i-1} = \mathbf{K} \mathbf{a}^{i-1} - \mathbf{f}$$

With (17.76) and (17.69), (17.77) therefore provides the following solution

$$\beta^{i-1} = -\frac{(\mathbf{s}^{i-1})^T \boldsymbol{\psi}^{i-1}}{(\mathbf{s}^{i-1})^T \mathbf{K} \mathbf{s}^{i-1}} \quad (17.78)$$

which is the best possible choice of  $\beta^{i-1}$ . If  $\boldsymbol{\psi}^{i-1} = \mathbf{0}$ , then the estimate  $\mathbf{a}^{i-1}$  fulfills the equilibrium equations and it is therefore the solution sought for. Accordingly, (17.78) provides  $\beta^{i-1} = 0$  when  $\boldsymbol{\psi}^{i-1} = \mathbf{0}$ . However, if  $\boldsymbol{\psi}^{i-1} \neq \mathbf{0}$  then (17.78) gives a solution for  $\beta^{i-1}$  which minimizes the potential energy  $\Pi$  as much as possible.

It appears that for linear problems, use of the search direction  $\mathbf{s}^{i-1}$  in the orthogonality condition is the best possible choice, cf. (17.77). It therefore seems



**Figure 17.17:** Line search for orthogonality between  $s^{i-1}$  and  $\psi(a^{i-1} + \beta^{i-1} s^{i-1})$ .

reasonable to adopt the same orthogonality condition also for other problems and this was already anticipated by the choice (17.70).

To further evaluate (17.70), we consider the scalar product  $(s^{i-1})^T \psi(a^{i-1})$  which with (17.67) can be written as

$$(s^{i-1})^T \psi(a^{i-1}) = (\psi(a^{i-1}))^T s^{i-1} = -(\psi(a^{i-1}))^T (A^{i-1})^{-1} \psi(a^{i-1}) \quad (17.79)$$

Since the iteration matrix  $A^{i-1}$  is usually a positive definite matrix, (17.79) shows that  $(s^{i-1})^T \psi(a^{i-1}) < 0$  usually holds. We then obtain the illustration shown in Fig. 17.17, where it is observed that  $\psi(a^{i-1})$  corresponds to  $\beta^{i-1} = 0$ . When  $\beta^{i-1}$  is changed, the out-of-balance forces change and  $\beta^{i-1}$  is changed such that  $s^{i-1}$  and  $\psi(a^{i-1} + \beta^{i-1} s^{i-1})$  eventually become orthogonal. Therefore, when  $\beta^{i-1}$  is determined from (17.70) this implies that the out-of-balance force has no components in the search direction  $s^{i-1}$ .

As already mentioned, (17.70) provides for, say, plasticity problems a non-linear relation for determination of  $\beta^{i-1}$ . Therefore, in principle, solution of (17.70) requires an iterative procedure and in each of these iterations, a new out-of-balance force vector needs to be determined. This in turn requires determination of the stress field, which means integration of the constitutive equations for all (Gauss) points in the body. In the next chapter, we shall treat integration of the nonlinear elasto-plastic constitutive equations in detail, but already at this stage it is evident that this integration is far from being trivial. Therefore, using a number of iterations to solve  $\beta^{i-1}$  from (17.70) in an accurate manner requires a computer effort that is not insignificant. Moreover, since the search direction  $s^{i-1}$  is only an approximation to the correct one, this also questions the meaningfulness of achieving an accurate  $\beta^{i-1}$ -solution. Finally, whereas the orthogonality expressed of (17.70) provides the best possible  $\beta^{i-1}$ -solution for

linear problems, it can only be expected to provide a fair  $\beta^{i-1}$ -solution for non-linear problems.

All these arguments suggest that there is no sense in trying to solve (17.70) in an accurate manner and, instead, an approximate value for  $\beta^{i-1}$  is sufficient. Indeed, this is supported by numerical experience both within structural mechanics, see Crisfield (1991), and optimization theory, see Luenberger (1984). This numerical experience also shows that also use of an approximate  $\beta^{i-1}$ -value, in general, reduces the number of equilibrium iterations significantly and reduces the total computer effort considerably. Due to its effectiveness, line search is available in most general purpose finite element programs.

With DOF being the number of degrees of freedom for the structure and B the bandwidth of the equation system, the computational cost for solving the equilibrium equation system by means of Gauss elimination is proportional to  $\text{DOF} \cdot B^2$ , cf. for instance Ottosen and Petersson (1992). The stresses  $\sigma^{i-1}$  are evaluated at each Gauss point in the structure and the number of Gauss points is roughly proportional to  $\text{DOF} = \text{degrees of freedom}$ . It is concluded that the advantage with the use of line search increases with the number of degrees of freedom.

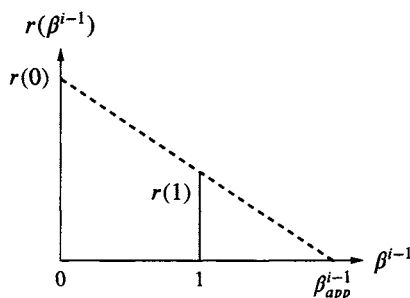


Figure 17.18: Approximate determination of  $\beta^{i-1}$ .

To arrive at an approximate  $\beta^{i-1}$ -value, we first observe that  $\beta^{i-1} = 1$  corresponds to the iteration scheme without line search, cf. (17.69) and (17.18). We therefore certainly expect that  $\beta^{i-1} > 0$ . Define the quantity  $r(\beta^{i-1})$  by

$$r(\beta^{i-1}) = -(s^{i-1})^T \psi(a^{i-1} + \beta^{i-1} s^{i-1})$$

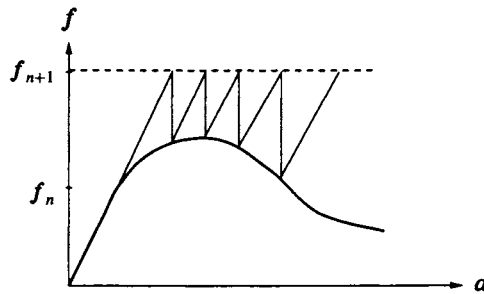
where the minus sign is motivated by the discussion following (17.79). The quantity  $\psi(a^{i-1})$  has already been calculated, cf. (17.67), and it is therefore easy to calculate  $r(\beta^{i-1} = 0)$ ; calculate also  $r$  for  $\beta^{i-1} = 1$ . We then obtain the situation sketched in Fig. 17.18. According to (17.70) we want  $r(\beta^{i-1}) = 0$  and making the linear extrapolation shown in Fig. 17.18, we obtain the following approximate value for  $\beta^{i-1}$

$$\beta_{app}^{i-1} = \frac{r(0)}{r(0) - r(1)} \quad (17.80)$$

If the difference  $r(0) - r(1)$  is small,  $\beta_{app}^{i-1}$  becomes large and, in principle, we may even encounter that  $\beta_{app}^{i-1}$  as predicted by (17.80) becomes negative. Since the entire motivation for the orthogonality condition (17.70) holds only strictly for linear problems, very small or very large  $\beta^{i-1}$ -values should be considered with caution. In practice therefore, acceptable values for  $\beta^{i-1}$  are restricted to a certain interval, say  $0.3 \leq \beta^{i-1} \leq 3$ . If (17.80) provides  $\beta_{app}^{i-1} < 0.3$  then  $\beta^{i-1} = 0.3$  is used and if (17.80) predicts  $\beta_{app}^{i-1} > 3$  then  $\beta^{i-1} = 3$  is used.

## 17.10 Limit points

In the previous discussion of incremental-iterative schemes, we have assumed that, for each load step, the load is increased by a certain amount and then held constant during the equilibrium iterations. This approach works well when the displacements increase with increasing load level.

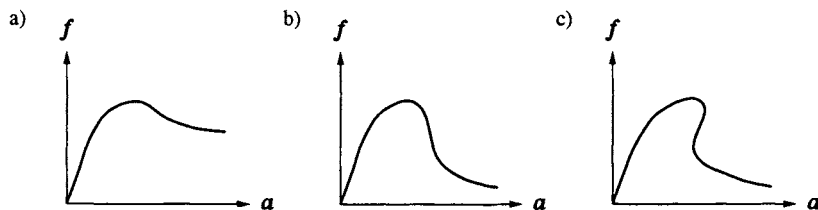


**Figure 17.19:** Structure displaying a softening response after the peak load. With the load  $f_{n+1}$  held fixed, equilibrium iterations will continue for ever (initial stiffness method used in the figure).

Consider now the situation where the structure displays a softening response, i.e. the structure has a maximum load capacity with a descending branch after the peak load, cf. Fig. 17.19. Evidently, the approach where the load is increased from  $f_n$  to  $f_{n+1}$  and then held constant will not be able to trace the softening branch; as soon as the load  $f_{n+1}$  exceeds the maximum load capacity of the structure, equilibrium iterations will simply continue for ever. This is illustrated in Fig. 17.19 using the initial stiffness method. Even if the response after the peak load does not display a descending branch, but exhibits a horizontal plateau (e.g. ideal plasticity), equilibrium iterations will continue for ever, if  $f_{n+1}$  exceeds the peak load.

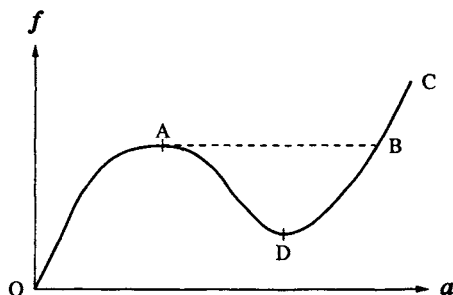
From an engineering point of view, the peak load, i.e. the *limit load*, is often of fundamental interest. Occasionally, it is argued that once the equilibrium iterations do not converge then the limit load has been reached. However, this

is a dangerous approach since equilibrium iterations may diverge for a number of other reasons. Moreover, it is often interesting to trace the descending branch after the peak load since it contains important information of the ductility or brittleness of the structural response, cf. Fig 17.20. In Fig. 17.20c) the phenomenon of *snap-back* is illustrated.



**Figure 17.20:** a) Ductile response; b) brittle response; c) snap-back.

The situation of *snap-through* often encountered in buckling problems is shown in Fig. 17.21. With the iterative approach discussed so far, one would only be able to trace the response OABC whereas the part ADB of the response cannot be traced.



**Figure 17.21:** Snap-through of structure

It is evidently of importance to derive iterative schemes that are able to trace the various post-peak responses discussed above. Strategies that possess this property are called *path-following methods* and borrowing from the terminology of buckling analysis they are also often called *continuation methods*.

Whereas we previously accepted the load level and then modified the displacements until equilibrium is fulfilled, it is evident that we must now devise iterative schemes where both the displacements and the load level are adjusted, i.e.

*Determination of the post-peak response requires simultaneous adjustment of the displacements and of the load level*

Before entering this discussion, it is of interest to be able to characterize the different kinds of response discussed in relation to Figs. 17.20 and 17.21. Assume for simplicity that the external loading is proportional, i.e.

$$\mathbf{f} = \lambda \mathbf{f}_{ref} \quad (17.81)$$

where  $\mathbf{f}_{ref}$  is a constant *reference load* and the proportionality factor  $\lambda = \lambda(t)$ , where  $t$  is the time, is called the *load parameter*. According to (17.2), the out-of-balance forces  $\boldsymbol{\psi}$  are then given by

$$\boldsymbol{\psi} = \mathbf{f}_{int} - \lambda \mathbf{f}_{ref}$$

where the internal forces  $\mathbf{f}_{int}$  depend on the stresses and thereby on the nodal displacements  $\mathbf{a}$ . The equilibrium condition is therefore given by

$$\boldsymbol{\psi} = \mathbf{0}; \quad \boldsymbol{\psi} = \boldsymbol{\psi}(\mathbf{a}, \lambda)$$

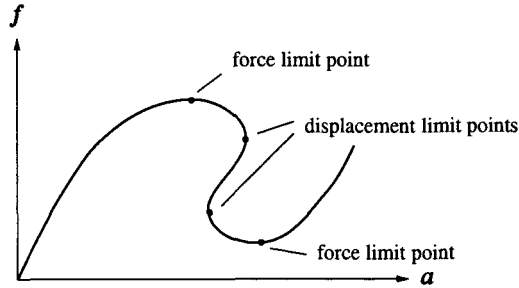
During equilibrium we have  $\boldsymbol{\psi} = \mathbf{0}$ , which leads to

$$\frac{\partial \boldsymbol{\psi}}{\partial \mathbf{a}} \dot{\mathbf{a}} + \frac{\partial \boldsymbol{\psi}}{\partial \lambda} \dot{\lambda} = \mathbf{0}$$

i.e.

$$\boxed{\mathbf{K}_t \dot{\mathbf{a}} - \dot{\lambda} \mathbf{f}_{ref} = \mathbf{0}} \quad (17.82)$$

where use was made of (17.31).



**Figure 17.22:** Illustration of force and displacement limit points

Let us make the following definitions:

$$\text{Force limit point: } \dot{\lambda} = 0 \text{ and } \dot{\mathbf{a}} \neq \mathbf{0}$$

$$\text{Displacement limit point: } \dot{\lambda} \neq 0 \text{ and } \dot{\mathbf{a}} = \mathbf{0}$$

(17.83)

These *limit points* are illustrated in Fig. 17.22 and we may note that in the literature such limit points are often called *turning points*.



To determine a force limit point, (17.82) and (17.83a) lead to  $\mathbf{K}_t \dot{\mathbf{a}} = \mathbf{0}$  and a non-trivial solution requires  $\det \mathbf{K}_t = 0$ , i.e.

$$\det \mathbf{K}_t = 0 \Rightarrow \text{force limit point}$$

Let us next assume that  $\det \mathbf{K}_t \neq 0$ ; then (17.82) implies

$$\dot{\mathbf{a}} = \lambda \mathbf{K}_t^{-1} \mathbf{f}_{ref}$$

A comparison with (17.83b) shows that the only possibility for a displacement limit point is that  $\mathbf{f}_{ref}$  fulfills  $\mathbf{K}_t^{-1} \mathbf{f}_{ref} = \mathbf{0}$ , which requires  $\det \mathbf{K}_t^{-1} = 0$ , i.e.

$$\det \mathbf{K}_t^{-1} = 0 \text{ and } \mathbf{K}_t^{-1} \mathbf{f}_{ref} = \mathbf{0} \Rightarrow \text{displacement limit point}$$

## 17.11 Bergan's minimum residual force method

In the previous iteration schemes, the load was increased step-wise and after each load increase, the load was held constant and, via equilibrium iterations, the nodal displacements were then adjusted in an iterative fashion until equilibrium was satisfied.

An elegant and entirely different approach, which allows force limit points to be passed as well as tracing of the post-peak response, was proposed by Bergan (1979, 1980). For a certain load level, the standard iteration scheme (17.18) is adopted, i.e.

$$\mathbf{A}^{i-1}(\mathbf{a}^i - \mathbf{a}^{i-1}) = -\boldsymbol{\psi}^{i-1} \quad (17.84)$$

where

$$\boldsymbol{\psi}^{i-1} = \mathbf{f}_{int}^{i-1} - \lambda^{i-1} \mathbf{f}_{ref}$$

and proportional loading according to (17.81) was assumed. It follows that

$$\boldsymbol{\psi}^0 = (\mathbf{f}_{int})_n - \lambda^0 \mathbf{f}_{ref} \quad (17.85)$$

where it was used that  $\mathbf{f}_{int}^0 = \int_V \mathbf{B}^T \boldsymbol{\sigma}^0 dV = \int_V \mathbf{B}^T \boldsymbol{\sigma}_n dV = (\mathbf{f}_{int})_n$ , cf. (17.22). Moreover, in (17.85) we set

$$\lambda^0 = \lambda_n + \Delta \lambda^* \quad (17.86)$$

where  $\mathbf{f}_n = \lambda_n \mathbf{f}_{ref}$  and the quantity  $\Delta \lambda^*$  is specified by us.

From (17.84), the displacements  $\mathbf{a}^i$  can now be determined; the corresponding stresses  $\boldsymbol{\sigma}^i$  can now also be obtained by integration of the constitutive equations. With the stress field  $\boldsymbol{\sigma}^i$  known throughout the body, the corresponding internal forces can be determined in the usual manner by means of

$$\mathbf{f}_{int}^i = \int_V \mathbf{B}^T \boldsymbol{\sigma}^i dV \quad (17.87)$$

The new residual forces then become

$$\psi^i = f_{int}^i - \lambda^i f_{ref} \quad (17.88)$$

where the new load parameter  $\lambda^i$  is viewed as unknown. The key point in *Bergan's minimum residual force method* is then to adjust the load parameter  $\lambda^i$  so that the residual force  $\psi^i$  becomes minimum.

Define the quantity  $b$  by

$$b = (\psi^i)^T \psi^i$$

With both  $f_{ref}$  and  $f_{int}^i$  being known and fixed and  $\lambda^i$  the only unknown, we obtain

$$\frac{db}{d\lambda^i} = -2(\psi^i)^T f_{ref}$$

as well as  $d^2b/(d\lambda^i)^2 = 2f_{ref}^T f_{ref} > 0$ . Therefore, the minimum of the quantity  $b$  with respect to  $\lambda^i$  is given by  $db/d\lambda^i = 0$ , i.e.

$$\boxed{(\psi^i)^T f_{ref} = 0} \quad (17.89)$$

Insertion of (17.88) gives the new load parameter  $\lambda^i$  according to

$$\boxed{\lambda^i = \frac{(f_{int}^i)^T f_{ref}}{f_{ref}^T f_{ref}}} \quad (17.90)$$

In view of the orthogonality expressed by (17.89) and the definitions of residual forces given by (17.88), we obtain the graphical illustration of the method shown in Fig. 17.23. After  $a^i$  has been determined from (17.84), these nodal displacements are accepted and the corresponding fixed internal forces  $f_{int}^i$  are then given by (17.87). The new load parameter  $\lambda^i$  is now adjusted according to (17.90) until the new external force  $\lambda^i f_{ref}$  and the new residual force  $\psi^i$  become orthogonal. Each iteration can therefore be viewed as a *two-step procedure*.

With this new load parameter, we perform a new iteration, i.e. the iteration number  $i$  is increased by one and the procedure is repeated until convergence – according to the procedures discussed in Section 17.6 – has been obtained. In that case, we put  $a_{n+1} = a^i$  and  $\lambda_{n+1} = \lambda^i$ .

Bergan's minimum residual method is simple and very effective. In (17.84), we can take any iteration matrix; for instance,  $A^{i-1}$  may be taken as the tangential stiffness matrix  $K_t^{i-1}$  or the constant initial stiffness matrix  $K$ . This latter choice is often adopted since it is cheap and since it allows force limit points to be passed and post-peak response to be traced. These features are illustrated in the one-dimensional case shown in Fig. 17.24.

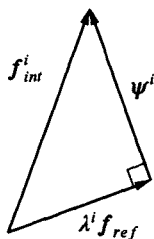


Figure 17.23: Orthogonal property obtained by Bergan's minimum residual method.

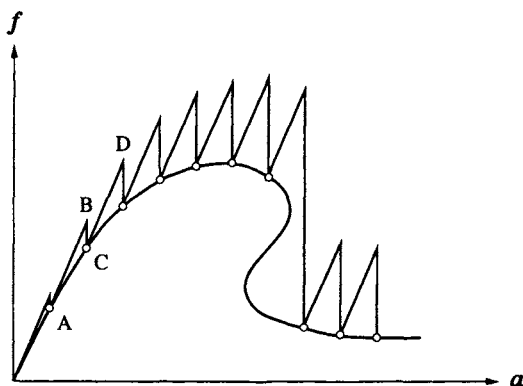


Figure 17.24: One-dimensional system. Bergan's method used with the initial stiffness matrix; after equilibrium is fulfilled, the imposed load increment is the same in all load steps. The predictions are indicated by (o).

Assume that point A has been reached. The load is then increased by some amount specified by us ( $\Delta\lambda^*$  in (17.86)) to point B and the corresponding displacement is determined by (17.84) using the initial stiffness. These displacements are accepted and the corresponding internal force is determined by (17.87); this internal force corresponds to point C. The load is now adjusted so that the residual force, i.e. BC, becomes minimum. In the present one-dimensional case, this means reducing the load from B to C, i.e. the new residual force becomes in the one-dimensional case exactly equal to zero. Therefore in one iteration, exact equilibrium has been obtained. The loading is now increased by an amount specified by us ( $\Delta\lambda^*$  in (17.86)) which brings us to point D and the procedure is repeated; in the case shown in Fig. 17.24 all specified load steps  $\Delta\lambda^*$  are taken as equal.

While exact fulfillment of equilibrium in just one iteration is not a general property, but holds only for the one-dimensional problems, Fig. 17.24 illustrates that Bergan's method allows us to pass the force limit point and to trace the post-peak response. However, the figure also illustrates that Bergan's method cannot trace a snap-back response.

---

**Box 17.3** Bergan's minimum residual force method
 

---

- *Initiation of quantities*

$$\mathbf{a}_0 = \mathbf{0}; \quad \boldsymbol{\epsilon}_0 = \mathbf{0}; \quad \boldsymbol{\sigma}_0 = \mathbf{0}; \quad \mathbf{f}_0 = \mathbf{0}; \quad \mathbf{f}_{int} = \mathbf{0}$$
  - *For load step  $n = 0, 1, 2, \dots, N_{max}$* 
    - *Determine new load level  $\lambda \mathbf{f}_{ref}$*
    - *Initiation of iteration quantities*

$$\mathbf{a}^0 := \mathbf{a}_n$$
    - *Iterate  $i = 1, 2, \dots$  until  $|\boldsymbol{\psi}|_{norm} = |\mathbf{f}_{int} - \lambda \mathbf{f}_{ref}|_{norm} < \epsilon_{\text{epsilon}}$* 
      - *Calculate  $\mathbf{K} = \int_V \mathbf{B}^T \mathbf{D} \mathbf{B} dV$*
      - *Calculate  $\mathbf{a}^i$  from  $\mathbf{K}(\mathbf{a}^i - \mathbf{a}^{i-1}) = \lambda \mathbf{f}_{ref} - \mathbf{f}_{int}$*
      - *Calculate  $\boldsymbol{\epsilon}^i := \mathbf{B} \mathbf{a}^i$*
      - *Determine  $\boldsymbol{\sigma}^i$  by integration of the constitutive equations (see next chapter)*
      - *Calculate internal forces  $\mathbf{f}_{int}$*
      - *Calculate load parameter  $\lambda = \frac{\mathbf{f}_{int}^T \mathbf{f}_{ref}}{\mathbf{f}_{ref}^T \mathbf{f}_{ref}}$*
    - *End iteration loop*
    - *Accept quantities*

$$\mathbf{a}_{n+1} := \mathbf{a}^i; \quad \boldsymbol{\epsilon}_{n+1} := \boldsymbol{\epsilon}^i; \quad \boldsymbol{\sigma}_{n+1} := \boldsymbol{\sigma}^i; \quad \mathbf{f}_{int}$$
  - *End load step loop*
- 

With the use of the initial stiffness matrix  $\mathbf{K}$ , each iteration using Bergan's method is in practice as cheap as iterations using the standard initial stiffness method (Bergan's method only requires the additional computation of  $\lambda^i$  given by (17.90), which is inexpensive). However, the standard initial stiffness method is not able to pass the force limit point and to trace the post-peak response. In addition to these features Bergan's method improves the convergence speed very significantly, as can be seen by comparing Figs. 17.10 and 17.24.

Here, we have formulated Bergan's method for proportional loading expressed by  $\mathbf{f} = \lambda \mathbf{f}_{ref}$ , but it is straightforward to generalize the method to non-proportional loading. In that case, the loading is considered as piece-wise linear and the above approach only requires modest modifications.

The algorithm for Bergan's method in combination with the initial stiffness matrix  $\mathbf{K}$  is summarized in Box. 17.3; here the means to enforce the kinematic

boundary conditions have not been indicated.

It is also noticed that a closely related method that improves the possibilities for passing displacement limit points was proposed by Krenk (1995) and Krenk and Hededal (1995). Instead of the orthogonality relation (17.89), this method makes use of the orthogonality relation  $(\psi^i)^T(\mathbf{a}^i - \mathbf{a}^{i-1}) = 0$  to determine the new load parameter  $\lambda^i$ .

Finally, we mention the very successful *arc-length* method where a combination of load and displacement control is adopted. This method facilitates the passage of limit points and a detailed discussion is provided, for instance by Crisfield (1991).