



# The Discrete Dynamical Functional Particle Method for Solving Constrained Optimization Problems

Mårten Gulliksson<sup>a</sup>

*Communicated by R. Cavoretto and M. Vianello*

## Abstract

The dynamical functional particle method (DFPM) is a method for solving equations by using a damped second order dynamical system. The dynamical system is solved by a symplectic method that is especially tailored for conservative systems. In this work we have extended DFPM to convex optimization problems with constraints. The method is tested on linear eigenvalue problems with normalization and orthogonality constraints as well as some simple nonlinear convex optimization problems.

## 1 Introduction

Consider the problem of solving a system of nonlinear equations

$$F(v) = 0, F : \mathbb{R}^n \rightarrow \mathbb{R}^n. \quad (1)$$

The idea behind the dynamical functional particle method (DFPM), [5], is to solve (1) through a damped second order dynamical system of the form

$$\ddot{u} + \eta \dot{u} = F(u) \quad (2)$$

such that  $\lim_{t \rightarrow t^*} F(u) = 0$ ,  $t^* \leq \infty$ . The parameter  $\eta > 0$  is the damping parameter that can be chosen in order to get fast convergence. It can be proven that if there is a convex potential  $V(u)$  such that  $F = -\nabla V$  then DFPM will converge to the minimum of  $V$ , i.e., the solution of

$$\min_{u \in \mathbb{R}^n} V(u),$$

or when  $\nabla V = -F = 0$ , see [5] and [7]. What makes DFPM using (2) especially appealing is the use of symplectic methods [6] since for a conservative system these methods will approximately retain the total energy of the system. This means that for a damped system a symplectic method will approximately follow the exact decay of the energy.

In [3] it is shown that DFPM using a symplectic Euler for finding the extreme states of linear and nonlinear Schrödinger equations is faster for very large problems than other competing methods. The results also shows that DFPM is as fast or faster than other methods for finding the smallest (or largest) eigenvalues of a linear eigenvalue problem. Further motivation for using DFPM is given in [2] where it is shown that the method is competitive with the best iterative methods for solving linear systems of equations. We stress that DFPM is very versatile and not restricted to linear problems. Indeed, in [8] we used DFPM to solve a highly nonlinear Schrödinger equation with nonlinear constraints and showed that it is a very interesting alternative to the so called imaginary time approach.

In this paper we will generalise DFPM to include constraints, i.e.,

$$\begin{aligned} \min_{u \in \mathbb{R}^n} V(u) \\ \text{s.t } g(u) = 0 \end{aligned} \quad (3)$$

where  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m < n$ . We will assume that (3) is convex. Problems of the form (3) have been solved in the context of the nonlinear Schrödinger describing Bose-Einstein condensate [8]. However, the approach in this paper is different since we do not use projection to satisfy the constraints but instead solve additionally damped dynamical equations, see next section. We will use symplectic Euler and show how the time step and damping parameter can be chosen in order to get an optimal local convergence rate. Based on our earlier successful results for non-constrained problems we strongly believe that DFPM for (3) will be competitive to other methods. However, in this paper we restrict ourselves to small problems to show the properties of the proposed method and leave a more extensive comparative test for future research.

The outline of the paper is the following. In Section 2 we derive DFPM for solving (3) and then use a linear analysis to give a necessary condition for asymptotic convergence. We then show how a Symplectic Euler can be used applied to a simple nonlinear problem with nonlinear constraints. In Section 3 we apply DFPM, the linear analysis, symplectic Euler and optimal parameters to the linear eigenvalue problem with norm and orthogonality constraints and give some numerical tests.

<sup>a</sup>School of Science and Technology, Örebro University, Sweden

## 2 DFPM using damped constraints

Define the Lagrange function of (3) as

$$L(u, \mu) = V(u) + g(u)^T \mu = V(u) + \sum_{j=1}^m g_j(u) \mu_j,$$

where  $\mu = (\mu_1(u), \dots, \mu_m(u))^T$  is a vector with the Lagrange parameters. The first order condition for a minimum is

$$\nabla_u L(u, \mu) = \nabla V(u) + \nabla g(u) \mu = 0, \quad g(u) = 0 \quad (4)$$

where  $\nabla g(u) = (\nabla g_1(u), \dots, \nabla g_m(u))$ . Here, we define gradients as column vectors and for  $m > 1$  the Jacobian  $\nabla g$  is the matrix with the columns  $\nabla g_j(u)$ . We will formulate DFPM for solving (3) based on the dynamical system

$$\ddot{u} + \eta \dot{u} = -\nabla V(u) - \nabla g(u) \mu \quad (5)$$

together with the constraints  $g(u) = 0$ . A well known approach to solve (5) is to use projection, i.e., at every time step  $k$  of the numerical scheme projection is used such that the iterate  $u_k$  implies  $g(u_k) = 0$ . For some simple constraints such as linear and quadratic constraints there are methods tailored for staying at the manifold but generally it is necessary to solve a nonlinear system of equations at each time step, for details see [6] Chapter IV.4 and references therein. We refer also to [8] where DFPM is applied to the nonlinear Schrödinger equation including quadratic constraints (normalization and angular momentum).

However, the method of projection is rather restrictive depending on the type of constraints and we would like to avoid the additional iterations for the projection. Therefore, we instead introduce the approach of satisfying the constraints by an additional damped dynamical system

$$\ddot{g}_i + \eta \dot{g}_i = -k_i g_i, \quad k_i > 0. \quad (6)$$

The equations in (6) are used to derive explicit expressions of  $\mu$  in (5), see the next section. In other words, the idea is to approach the manifold dynamically and not enforcing the constraints at every time step.

We use (6) to find explicit expressions for the Lagrange parameters  $\mu$  introduced in the introduction. We have

$$\dot{g}_i = \nabla g_i^T \dot{u}, \quad \ddot{g}_i = \dot{u}^T \nabla^2 g_i \dot{u} + \nabla g_i^T \ddot{u}$$

and inserting this and the expression for  $\ddot{u}$  in (5) in (6) gives

$$\dot{u}^T \nabla^2 g_i \dot{u} + \nabla g_i^T (-\eta \dot{u} - \nabla V - \nabla g_i \mu) + \eta \nabla g_i^T \dot{u} = -k_i g_i$$

or

$$\nabla g_i^T (u) \nabla g_i(u) \mu = k_i g_i(u) - \nabla g_i(u)^T \nabla V(u) + \dot{u}^T \nabla^2 g_i(u) \dot{u}, \quad i = 1, \dots, m. \quad (7)$$

The equations in (7) can be gathered together into a system

$$\nabla g(u)^T \nabla g(u) \mu(u, \dot{u}) = K g(u) - \nabla g(u)^T \nabla V(u) + h(u, \dot{u}), \quad (8)$$

$$K = \text{diag}(k_1, \dots, k_n), \quad h_i(u, \dot{u}) = \dot{u}^T \nabla^2 g_i \dot{u}$$

where a unique solution exists if  $\nabla g$  has full rank. Solving for  $\mu$  gives

$$\mu(u, \dot{u}) = -(\nabla g(u)^T \nabla g(u))^{-1} K g(u) + (\nabla g(u)^T \nabla g(u))^{-1} \nabla g(u)^T \nabla V(u) + (\nabla g(u)^T \nabla g(u))^{-1} h(u, \dot{u}), \quad (9)$$

where the second term is  $\nabla g(u)^+ \nabla V(u)$ . If we assume that  $u \rightarrow u^*$ ,  $t \rightarrow \infty$  the two last terms tends to zero, i.e.,  $\mu(u, \dot{u}) \rightarrow \nabla g(u^*)^+ \nabla V(u^*)$ . Furthermore, we note that the last term tends to zero to second order in  $\dot{u}$ . Moreover,  $\nabla V(u) + \nabla g(u) \mu \rightarrow (I - \nabla g(u^*) \nabla g(u^*)^+) \nabla V(u^*)$ , i.e., the projected gradient.

### 2.1 Linear stability analysis

Let us formulate DFPM for our problem as

$$\ddot{u} + \eta \dot{u} = -\nabla V - \nabla g \mu(u, \dot{u}). \quad (10)$$

where  $\mu(u, \dot{u})$  is given by the solution of the system in (8). We can write (10) as the system

$$\dot{u} = v \quad (11)$$

$$\dot{v} = -\nabla V - \nabla g \mu(u, v) - \eta v$$

with corresponding Jacobian

$$J(u, v) = \begin{pmatrix} 0 & I \\ -\nabla^2 V(u) - \sum_{j=1}^n \nabla^2 g_j(u) \mu_j(u, v) - \nabla g(u) \nabla_u \mu(u, v)^T & -v \nabla_v \mu(u, v)^T - \eta I \end{pmatrix}. \quad (12)$$

At the stationary point we have

$$\mu(u^*, 0) = (\nabla g^T \nabla g)^{-1} \nabla g^T \nabla V,$$

giving

$$J(u^*, 0) = \begin{pmatrix} 0 & I \\ -\nabla^2 V(u^*) - \sum_{j=1}^n \nabla^2 g_j(u^*) \mu_j(u^*, 0) - \nabla g \nabla_u \mu(u^*, 0)^T & -\eta I \end{pmatrix}$$

and we conclude that the method is asymptotically stable if the eigenvalues of

$$B = \nabla^2 V(u^*) + \sum_{j=1}^n \nabla^2 g_j(u^*) \mu_j(u^*, 0) + \nabla g(u^*) \nabla_u \mu(u^*, 0)^T \quad (13)$$

are all positive.

## 2.2 Symplectic Euler

What makes DFPM using (5),(6) especially appealing is the use of symplectic methods [6] since for a conservative system these methods will approximately retain the total energy of the system. This means that for a damped system a symplectic method will approximately follow the exact decay of the energy (decay because of dissipation).

It is straight forward to apply a symplectic Euler [6] method to the system (11). We assume that we have a constant time step  $\Delta t$  and for simplicity a constant damping parameter  $\eta$ . Using (11) the symplectic Euler reads

$$\begin{aligned} u_{k+1} &= u_k + \Delta t v_k \\ v_{k+1} &= v_k + \Delta t (-\nabla V(u_{k+1}) - \nabla g(u_{k+1})\mu(u_{k+1}, v_k) - \eta v_k) \end{aligned} \quad (14)$$

Close to the stationary solution  $u = u^*, v^* = 0$  the convergenc of (14) will be determined by the linear system

$$\dot{w} + \eta \dot{w} = -Bw$$

with  $B$  given in (13). Assume that the matrix  $B$  has eigenvalues  $\lambda_i(B) > 0$  and the dynamical system is asymptotically convergent. From a linear analysis of the Symplectic Euler method, see [2], we get that the number of iterations until convergence is

$$n \propto \sqrt{\frac{\lambda_{\max}(B)}{\lambda_{\min}(B)}}. \quad (15)$$

if parameters are chosen optimally as

$$\eta_{opt} = \frac{2\sqrt{\lambda_{\min}(B)\lambda_{\max}(B)}}{\sqrt{\lambda_{\min}(B)} + \sqrt{\lambda_{\max}(B)}}, \Delta t_{opt} = \frac{2}{\sqrt{\lambda_{\min}(B)} + \sqrt{\lambda_{\max}(B)}}. \quad (16)$$

It is not possible to find parameters optimally as in (16) for the general formulation of DFPM in (5) since the eigenvalues of  $B$  are only accessible for specific problems settings. However, it should be possible to iteratively choose the damping and time step according to some merit function used in Lagrange Multiplier Methods, see [1] but this will be future research. For the linear eigenvalue problem we discuss the problem of optimal parameter choice further in Section 3.3.

## 2.3 Nonlinear optimization example

We consider the simple nonlinear problem

$$\begin{aligned} \min_{u \in \mathbb{R}^4} V(u) &= e^{u_1^2 + u_2^2 + u_3^2 + u_4^2} \\ \text{s.t } g_1 &= u_3 - e^{-u_1 - u_2} = 0, g_2 = u_4 - e^{-u_2 - u_3} = 0. \end{aligned} \quad (17)$$

The solution  $u^* = (0.21859539, 0.37508690, 0.55228984, 0.39559008)^T$  is found both by using DFPM and Matlabs function *fmincon*. The eigenvalues of  $B$  in (13) where calculated using the numerical solution and the parameters was then chosen as  $k_1 = k_2 = 1, \eta = \eta_{opt}$  where  $\Delta t = \Delta t_{opt}$  where  $\eta_{opt}, \Delta t_{opt}$  are given by (29). The initial guess in DFPM was chosen with uniformly distributed elements in  $[0, 1]$ . The convergence is shown in Figure 1. We note that convergence does not imply that the norm of

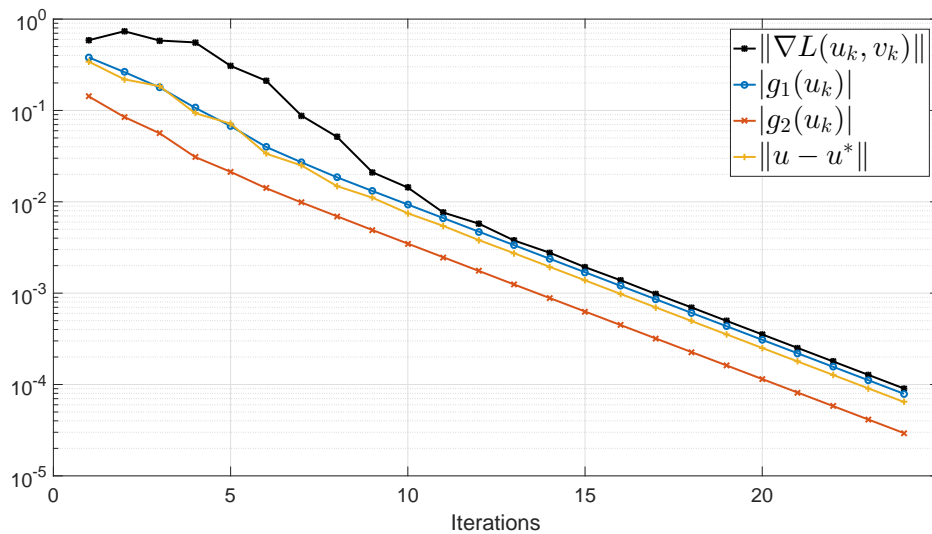


Figure 1: Convergence for problem (17) with  $\eta$  and  $\Delta t$  chosen optimally according to (29).

the Lagrangian is monotonically decreasing. Indeed, the global convergence properties of DFPM using damped constraints shown here are not known.

### 3 Non-defect linear eigenvalue problems

Consider the eigenvalue problem

$$Au = \lambda u, u^T u = 1, \quad (18)$$

where  $A \in \mathbb{R}^{n \times n}$  is assumed to be positive definite. Assume that the eigenvalues are distinct (non-defect eigenvalue problem) and sorted as  $0 < \lambda_1 < \lambda_2 < \dots < \lambda_n$  then the solution of

$$\begin{aligned} \min_{u \in \mathbb{R}^n} V(u) &= \frac{1}{2} u^T A u \\ \text{s.t } g(u) &= \frac{1}{2} (u^T u - 1) = 0 \end{aligned} \quad (19)$$

is the eigenvector  $u_1$  corresponding to the smallest eigenvalue  $\lambda_1$ . Furthermore,

$$\lambda_1 = u_1^T A u_1.$$

#### 3.1 First eigenvalue using orthogonality constraint

We follow the general approach and from (18) we define the corresponding Lagrange function

$$L(u, \mu) = \frac{1}{2} u^T A u + \mu \frac{1}{2} (u^T u - 1).$$

where

$$\nabla L_u(u, \mu) = Au + \mu u = 0.$$

Using the idea of damped constraints we formulate DFPM as

$$\ddot{u} + \eta \dot{u} = -Au - \mu u \quad (20)$$

$$\ddot{g} + \eta \dot{g} = -kg, \quad (21)$$

where  $k > 0$ . We have

$$\nabla g = u, \nabla^2 g = I, \nabla V = Au$$

that inserted in (7) gives

$$u^T u \mu = k \frac{1}{2} (u^T u - 1) - u^T A u + \dot{u}^T \dot{u}$$

or

$$\mu(u, \dot{u}) = \frac{1}{u^T u} (-u^T A u + k \frac{1}{2} (u^T u - 1) + \dot{u}^T \dot{u}).$$

Note that  $\mu(u) \rightarrow -\lambda_1$  if  $g(u) \rightarrow 0, \dot{u} \rightarrow 0$ . From (20) we get

$$\ddot{u} + \eta \dot{u} = -Au - \mu(u, \dot{u})u$$

that we may write as a system

$$\dot{u} = v$$

$$\dot{v} = -Au - \mu(u, v)u - \eta v$$

(22)

At a stationary point we have  $u = u_1, \dot{u} = v = 0, \mu = -\lambda_1$  which is, of course, the solution of the original eigenvalue problem.

**Theorem 3.1.** DFPM given by (20), (21) is asymptotically convergent to the eigenvector  $u_1$  and  $\mu_1(u_1, 0) = -\lambda_1$  in (9). Furthermore, the matrix  $B$  in (13) has positive eigenvalues  $k, \lambda_2 - \lambda_1, \lambda_3 - \lambda_1, \dots, \lambda_n - \lambda_1$ .

*Proof.* We need to show that (22) is stable at  $u = u_1, v = 0$ . Therefore, we consider the Jacobian of the right hand side in (22)

$$J(u, v) = \begin{pmatrix} 0 & I \\ -A - \mu I - u \nabla_u \mu^T & -\eta I - v \nabla_v \mu^T \end{pmatrix}$$

where

$$\nabla_u \mu = \frac{1}{u^T u} (-2Au + ku - 2u\mu), \nabla_v \mu = \frac{1}{u^T u} (2v)$$

and

$$\nabla_u \mu(u_1) = k u_1, \nabla_v \mu(u_1) = 0$$

which gives

$$J(u_1, 0) = \begin{pmatrix} 0 & I \\ -A + \lambda_1 I - k u_1 u_1^T & -\eta I \end{pmatrix} \quad (23)$$

Since the damped linear system

$$\ddot{u} + \eta \dot{u} = (-A + \lambda_1 I - (k + \lambda_1) u_1 u_1^T) u$$

has the corresponding Jacobian as in (23) DFPM will converge locally if the eigenvalues of

$$B = A - \lambda_1 I + k u_1 u_1^T$$

are positive. Let  $A = U \Lambda U^T$  be a similarity transformation of  $A$  where  $U^T U = I$  and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ . Then we have

$$U^T B U = \text{diag}(\lambda_1, \dots, \lambda_n) - \lambda_1 I + k U^T u_1 (U^T u_1)^T$$

and since  $U^T u_1 = e_1$

$$U^T B U = \text{diag}(k, \lambda_2 - \lambda_1, \dots, \lambda_{n-1} - \lambda_1, \lambda_n - \lambda_1)$$

i.e., we have derived the eigenvalues of  $B$  where  $\lambda(B) > 0$  and that DFPM is locally asymptotic convergent.  $\square$

### 3.2 Finding the whole spectrum

Assume that we know the eigenvalues  $\lambda_1, \dots, \lambda_{m-1}$  and the normalized eigenvectors  $u_1, \dots, u_{m-1}$  and want to find  $\lambda_m$  and corresponding eigenvector  $u_m$ . The unknown eigenvalue and eigenvector is given by the minimization problem

$$\begin{aligned} \min_{u \in \mathbb{R}^n} V(u) &= \frac{1}{2} u^T A u \\ \text{s.t } u_1^T u &= 0, \dots, u_{m-1}^T u = 0, (u^T u - 1)/2 = 0. \end{aligned} \quad (24)$$

From (5) we get DFPM as

$$\ddot{u} + \eta \dot{u} = -Au + \sum_{i=1}^{m-1} \mu_i u_i + \mu_m u \quad (25)$$

$$\ddot{g}_i + \eta \dot{g}_i = -k_i g_i, i = 0, \dots, m-1 \quad (26)$$

where  $g_m = (u^T u - 1)/2 = 0$ ,  $g_i = u_i^T u = 0, i = 1, \dots, m-1$  giving  $\nabla g_m = u$ ,  $\nabla g_i = u_i, i = 1, \dots, m-1$ .

**Theorem 3.2.** DFPM given by (25), (26) is asymptotically convergent to the eigenvector  $u_m$  and  $\mu_m(u_m, 0) = -\lambda_m$  in (9). Furthermore, the matrix  $B$  in (13) has positive eigenvalues  $k_1, \dots, k_{m-1}, k_m, \lambda_{m+1} - \lambda_m, \lambda_{m+2} - \lambda_m, \dots, \lambda_n - \lambda_m$ .

*Proof.* We want to show that DFPM in this case is asymptotically stable and therefore we need to find the Jacobian

$$J(u, v) = \begin{pmatrix} 0 & I \\ -\nabla^2 V - \sum_{j=1}^n \nabla^2 g_j \mu_j(u, v) - \nabla g \nabla_u \mu^T & -v \nabla_v \mu^T - \eta I \end{pmatrix}.$$

and evaluate it at the stationary point  $u = u_m, v = 0$ . The difficult part is to find the derivative  $\nabla_u \mu$  and thus we first consider the expression for  $\mu$ . It is easy to see that

$$\nabla g = (u_1, \dots, u_{m-1}, u) = (U, u) \in \mathbb{R}^{n \times m}$$

with an obvious definition of  $U$  giving

$$\nabla g^T \nabla g = \begin{pmatrix} U^T \\ u^T \end{pmatrix} (U, u) = \begin{pmatrix} I & U^T u \\ u^T U & u^T u \end{pmatrix}$$

with an inverse

$$(\nabla g^T \nabla g)^{-1} = \frac{1}{u^T u - u^T U U^T u} \begin{pmatrix} u^T u I + U^T u u^T U & -U^T u \\ -u^T U & 1 \end{pmatrix}.$$

Furthermore, from (9) we have

$$\mu = -(\nabla g^T \nabla g)^{-1} K g + (\nabla g^T \nabla g)^{-1} \nabla g^T \nabla V + (\nabla g^T \nabla g)^{-1} h,$$

where

$$K = \text{diag}(k_1, \dots, k_m), h_i = -\dot{u}^T \nabla^2 g_i \dot{u}$$

which in this case looks like

$$\mu = \frac{1}{u^T u - u^T U U^T u} \begin{pmatrix} u^T u I_{m-1} + U^T u u^T U & -U^T u \\ -u^T U & 1 \end{pmatrix} \left( -K g + \begin{pmatrix} U^T \\ u^T \end{pmatrix} A u + h(u, \dot{u}) \right),$$

$$h_i(u, \dot{u}) = 0, i = 1, \dots, m-1, h_m = -\dot{u}^T \dot{u}.$$

We now derive  $\nabla_u \mu(u_m, 0)$ . At the stationary point we have  $g(u_m) = 0, U^T u_m = 0, U^T A u = \Lambda_{m-1} U^T u_m = 0$  and taking this into account by dropping all terms zero at  $u_m$  we do an implicit derivation and get

$$u^T u \nabla_u \mu + 2\mu u^T = -K \nabla g^T + \frac{\partial}{\partial u} \begin{pmatrix} u^T u I & -U^T u \\ -u^T U & 1 \end{pmatrix} \begin{pmatrix} U^T \\ u^T \end{pmatrix} A u$$

or

$$\begin{aligned} u^T u \nabla_u \mu + 2\mu u^T &= -K \begin{pmatrix} U^T \\ u^T \end{pmatrix} + \frac{\partial}{\partial u} \begin{pmatrix} u^T u I & -U^T u \\ -u^T U & 1 \end{pmatrix} \begin{pmatrix} U^T A u \\ u^T A u \end{pmatrix} = \\ &= -K \begin{pmatrix} U^T \\ u^T \end{pmatrix} + \frac{\partial}{\partial u} \begin{pmatrix} u^T u U^T A u - u^T A u U^T u \\ -u^T U U^T A u + u^T A u \end{pmatrix} = \\ &= -K \begin{pmatrix} U^T \\ u^T \end{pmatrix} + \begin{pmatrix} u^T u U^T A - u^T A u U^T \\ 2u^T A \end{pmatrix} \end{aligned}$$

At the stationary point we have  $u^T A u = \lambda_m, U^T A = \text{diag}(\lambda_1, \dots, \lambda_{m-1}) U^T, A u = \lambda_m u_m, u^T u = 1$  giving

$$u^T u \nabla_u \mu + 2\mu u^T = -K \begin{pmatrix} U^T \\ u^T \end{pmatrix} + \begin{pmatrix} u^T u U^T A - u^T A u U^T \\ 2u^T A \end{pmatrix}$$

and inserting  $u = u_m$

$$\nabla_u \mu(u_m, 0) = -K \begin{pmatrix} U^T \\ u_m^T \end{pmatrix} + \begin{pmatrix} \Lambda_{m-1} U^T - \lambda_m U^T \\ 2\lambda_m u_m^T \end{pmatrix} - 2\mu(u_m, 0) u_m^T.$$

Since  $\nabla g(u_m) = (U, u_m)$  we get

$$\nabla g(u_m) \nabla_u \mu(u_m, 0)^T = -(U, u_m) K \begin{pmatrix} U^T \\ u_m^T \end{pmatrix} + (U, u_m) \begin{pmatrix} \Lambda_{m-1} U^T - \lambda_m U^T \\ 2\lambda_m u_m^T \end{pmatrix} - 2(U, u_m) \mu(u_m, 0) u_m^T$$

or with  $\mu(u_m, 0)^T = (0_{m-1}, \lambda_m)$

$$\begin{aligned} \nabla g(u_m) \nabla_u \mu(u_m, 0)^T &= -(U, u_m) K \begin{pmatrix} U^T \\ u_m^T \end{pmatrix} + U \Lambda_{m-1} U^T - \lambda_m U U^T + 2\lambda_m u_m u_m^T - 2\lambda_m u_m u_m^T = \\ &= -(U, u_m) K \begin{pmatrix} U^T \\ u_m^T \end{pmatrix} + U \Lambda_{m-1} U^T - \lambda_m U U^T. \end{aligned}$$

Furthermore,

$$\sum_{j=1}^n \nabla^2 g_j(u_m) \mu_j(u_m, 0) = \nabla^2 g_m(u_m) \mu_m(u_m, 0) = -\lambda_m I$$

giving the matrix of interest

$$B = A + (U, u_m) K \begin{pmatrix} U^T \\ u_m^T \end{pmatrix} - U \Lambda_{m-1} U^T + \lambda_m U U^T + \lambda_m I.$$

A similarity transformation with  $Q = (U, \bar{U})$  gives

$$\begin{aligned} Q^T B Q &= \Lambda + \begin{pmatrix} U^T \\ \bar{U}^T \end{pmatrix} (U, u_m) K \begin{pmatrix} U^T \\ u_m^T \end{pmatrix} (U, \bar{U}) - \begin{pmatrix} U^T \\ \bar{U}^T \end{pmatrix} U \Lambda_{m-1} U^T (U, \bar{U}) + \\ &+ \lambda_m \begin{pmatrix} U^T \\ \bar{U}^T \end{pmatrix} U U^T (U, \bar{U}) - \lambda_m I. \end{aligned}$$

or

$$\begin{aligned} Q^T B Q &= \Lambda + \begin{pmatrix} I_{m-1} & 0 \\ 0 & \bar{U}^T u_m \end{pmatrix} K \begin{pmatrix} I_{m-1} & U^T \\ 0 & u_m^T \bar{U} \end{pmatrix} - \begin{pmatrix} I_{m-1} \\ 0 \end{pmatrix} \Lambda_{m-1} (I_{m-1}, 0) + \\ &+ \lambda_m \begin{pmatrix} I_{m-1} \\ 0 \end{pmatrix} (I_{m-1}, 0) - \lambda_m I. \end{aligned}$$

Using  $\bar{U}^T u_m = e_1$  we get

$$\begin{aligned} Q^T B Q &= \Lambda + \begin{pmatrix} I_{m-1} & 0 \\ 0 & e_1 \end{pmatrix} \begin{pmatrix} K_{m-1} & 0 \\ 0 & k_m \end{pmatrix} \begin{pmatrix} I_{m-1} & U^T \\ 0 & e_1^T \end{pmatrix} - \begin{pmatrix} \Lambda_{m-1} & 0 \\ 0 & 0 \end{pmatrix} + \\ &+ \lambda_m \begin{pmatrix} I_{m-1} & 0 \\ 0 & 0 \end{pmatrix} - \lambda_m I \end{aligned}$$

or

$$Q^T B Q = \Lambda + \begin{pmatrix} K_{m-1} & 0 \\ 0 & k_m e_1 e_1^T \end{pmatrix} - \begin{pmatrix} \Lambda_{m-1} & 0 \\ 0 & 0 \end{pmatrix} + \lambda_m \begin{pmatrix} I_{m-1} & 0 \\ 0 & 0 \end{pmatrix} - \lambda_m I$$

which can be simplified to

$$Q^T B Q = \text{diag}(k_1, \dots, k_{m-1}, k_m, \lambda_{m+1} - \lambda_m, \lambda_{m+2} - \lambda_m, \dots, \lambda_n - \lambda_m) \quad (27)$$

showing that  $\lambda(B) > 0$  and that DFPM is locally asymptotic convergent and the eigenvalues of  $B$ .  $\square$

### 3.3 Optimal parameter choice

Consider the matrix in (27), i.e., the eigenvalues of  $B$ . From (15) we see that in order not to slow down the convergence the damping parameters  $k_i$  should satisfy

$$\lambda_{\min}(B) < k_i < \lambda_{\max}(B), i = 1, \dots, m$$

and with this choice the optimal parameters are determined by the remaining  $\lambda_i(B) = \lambda_i - \lambda_m, i = m+1, \dots, n$ . Again, from the linear analysis we get

$$n \propto \sqrt{\frac{\lambda_n - \lambda_m}{\lambda_{m+1} - \lambda_m}} \quad (28)$$

and the optimal parameter choice is

$$\eta_{\text{opt}} = \frac{2\sqrt{(\lambda_{m+1} - \lambda_m)(\lambda_n - \lambda_m)}}{\sqrt{\lambda_{m+1} - \lambda_m} + \sqrt{\lambda_n - \lambda_m}}, \Delta t_{\text{opt}} = \frac{2}{\sqrt{\lambda_{m+1} - \lambda_m} + \sqrt{\lambda_n - \lambda_m}}. \quad (29)$$

Calculating the optimal parameters require the eigenvalues  $\lambda_{m+1}, \lambda_m, \lambda_n$  of  $A$  which obviously are not possible to find exactly in general. However, in many applications  $A$  is attained from a discretization procedure and the size is depending on some parameter  $h$  where smaller  $h$  will give a larger size of  $A$ . Therefore, an approximation of the eigenvalues can be efficiently extrapolated for some larger  $h$  and used to get approximate damping and time step parameters, see [2]. If extrapolation is not applicable and  $A$  has known properties such as diagonally dominance, perturbation analysis is an option [2].

### 3.4 Numerical tests

We have chosen to show how DFPM works for the linear eigenvalue problem with norm constraint and orthogonality constraint. The purpose of these tests is not to show the efficiency of the method, which can be seen in [8, 2], only to give some examples of the general approach using damped constraint equations. The parameters in the method as well as the time step of the symplectic method are chosen more or less optimal and it is left for future research to derive efficient ways of adaptively calculate these parameters during the iterations. The problem we solve numerically is

$$\begin{aligned} \min_{u \in \mathbb{R}^n} V(u) &= \frac{1}{2} u^T A u \\ \text{s.t } g_1 &= u_1^T u = 0, g_2 = (u^T u - 1)/2 = 0. \end{aligned} \quad (30)$$

where  $A \in \mathbb{R}^{100 \times 100}$  is randomly chosen as symmetric positive definite and the eigenvector  $u_1$  corresponding to the smallest eigenvalue is numerically calculated in Matlab for convenience (could of course be found by DFPM). The parameters  $k_1 = k_2 = 2.4409$ ,  $\eta = \eta_{opt}$  and  $\Delta t = \Delta t_{opt}$  where  $\eta_{opt}, \Delta t_{opt}$  are given by (29). The convergence is shown in Figure 2.

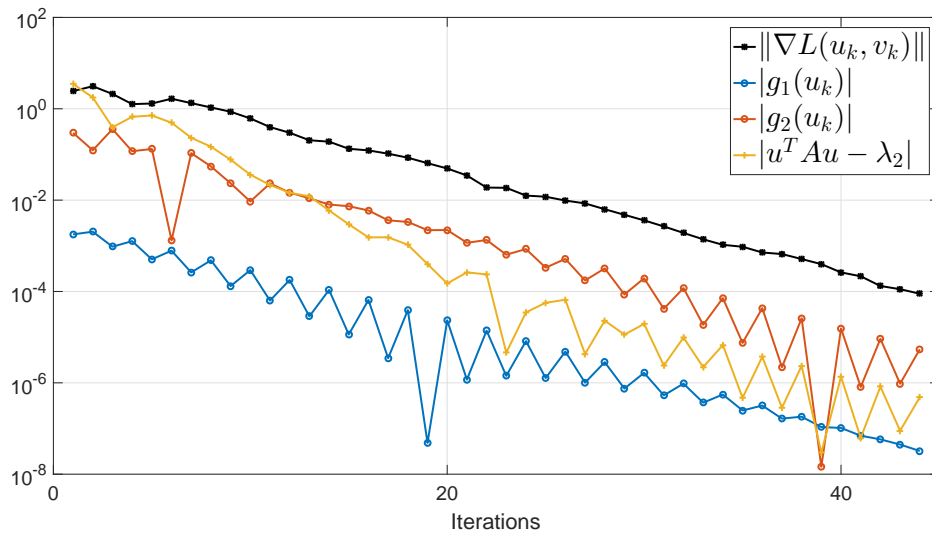


Figure 2: Convergence for problem (30) with  $\eta$  and  $\Delta t$  chosen optimally according to (29).

We remark that there are starting points that do not give convergence. Further, we see that the error in the eigenvalue is very much more accurate than the norm of the Lagrangian which is consistent with the behaviour of other methods for finding eigenvalues, see [4].

### References

- [1] D.P Bertsekas. *Constrained Optimization and Lagrange Multiplier Methods*. Athena scientific series in optimization and neural computation. Athena Scientific, 1996.
- [2] S. Edvardsson, M. Neuman, P. Edström, and H. Olin. Solving equations through particle dynamics. *Computer Physics Communications*, 197:169 – 181, 2015.
- [3] Sverker Edvardsson, Mårten Gulliksson, and Johan Persson. The dynamical functional particle method: An approach for boundary value problems. *Journal of applied mechanics*, 79(2):021012, 2012.
- [4] Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, 4th edition, 2013.
- [5] Mårten Gulliksson, Sverker Edvardsson, and Andreas Lind. The Dynamical Functional Particle Method. *ArXiv e-prints*, March 2013.
- [6] E. Hairer, C. Lubich, and G. Wanner. *Geometric Numerical Integration, 2nd ed*. Springer, 2006.
- [7] A. N. Michel, L. Hou, and D. Liu. *Stability of Dynamical Systems*. Birkhauser, 2008.
- [8] Patrik Sandin, Magnus Ögren, and Mårten Gulliksson. Numerical solution of the stationary multicomponent nonlinear schrödinger equation with a constraint on the angular momentum. *Phys. Rev. E*, 93:033301, March 2016.