

# Numerical study of liquid crystal elastomers by a mixed finite element method

C. LUO<sup>1</sup> and M. C. CALDERER<sup>2</sup>

<sup>1</sup> *Oxford Center for Collaborative Applied Mathematics, Mathematical Institute, University of Oxford, Oxford OX1 3LB, UK*

*email: luo@maths.ox.ac.uk*

<sup>2</sup> *School of Mathematics, University of Minnesota, Minneapolis, MN 55455, USA*

(Received 14 July 2011; revised 18 July 2011; accepted 19 July 2011;  
first published online 22 August 2011)

Liquid crystal elastomers present features not found in ordinary elastic materials, such as semi-soft elasticity and the related stripe domain phenomenon. In this paper, the two-dimensional Bladon–Terentjev–Warner model and the one-constant Oseen–Frank energy expression are combined to study the liquid crystal elastomer. We also impose two material constraints, the incompressibility of the elastomer and the unit director norm of the liquid crystal. We prove existence of minimiser of the energy for the proposed model. Next we formulate the discrete model, and also prove that it possesses a minimiser of the energy. The inf-sup values of the discrete linearised system are then related to the smallest singular values of certain matrices. Next the existence and uniqueness of the Lagrange multipliers associated with the two material constraints are proved under the assumption that the inf-sup conditions hold. Finally numerical simulations of the clamped-pulling experiment are presented for elastomer samples with aspect ratio 1 or 3. The semi-soft elasticity is successfully recovered in both cases. The stripe domain phenomenon, however, is not observed, which might be due to the relative coarse mesh employed in the numerical experiment. Possible improvements are discussed that might lead to the recovery of the stripe domain phenomenon.

**Key words:** Liquid crystal elastomer; Semi-soft elasticity; Variational methods; Mixed finite element method; Inf-sup condition.

## 1 Introduction

Liquid crystal elastomer (LCE) is an elastic material containing nematic liquid crystal molecules. Rotation of these liquid crystal molecules may lead to unique mechanical, optical and electrical properties, which are not observed in ordinary elastic materials. By properly exploiting these special properties, one might be able to manufacture new devices, such as artificial muscles [23].

The stripe domain formation and the semi-soft elasticity property are two special phenomena exhibited by LCEs [34]. They have been observed in the *clamped-pulling* experiment, in which a piece of rectangular LCE is clamped and pulled in the direction perpendicular to the initial uniform orientation of the liquid crystal directors. The stripe domain phenomenon refers to the formation of stripes with alternating director angles during the pulling process. In each stripe, the directors align along the same direction, while the directors in adjacent stripes align symmetrically about their middle line. For

square-shaped polysiloxane LCE, it is observed [25, 35] that during the pulling process, the stripe domain first occurred in the centre of the domain and then broke into two, which then migrated towards the two clamped ends. The semi-soft elasticity refers to the unusual stress–strain relationship during the pulling process. The LCE is first *hard*, in a regime such that the stress grows almost linearly with strain; then it reaches the *soft* regime, in which the stress remains almost constant while the strain increases; then the LCE becomes hard again upon further increase of the strain [12, 26].

Several models [3, 6, 7, 18, 33] have been proposed to explain the special behaviours of LCE. A very successful one is that proposed by Bladon, Terentjev and Warner (BTW) [3] that predicts the stripe domain phenomenon and the soft elastic response of the elastomer. However, the stress–strain relationship computed with the BTW model is *ideally soft* [34]. That is it lacks the initial hard regime typically observed in experiments. Several approaches have been proposed to modify the BTW and fully capture the experimental results. The Verwey–Warner–Terentjev model [33] extends the BTW model by adding a term related to the cross-linking state, and successfully recovers the semi-soft phenomenon. Other models [6, 7, 18] extend the BTW model by adding liquid crystal elastic energy terms, such as Oseen–Frank [21], Ericksen [19] and Landau–de Gennes energy terms [15]. Unlike the Verwey–Warner–Terentjev model, the latter ones involve first derivatives of the director field.

There are, however, relatively few works in the literature about the numerical simulation of LCEs. This may be due to the complexity caused by the coupling between the displacement of the bulk and the orientation of the liquid crystal directors. An important work on numerical simulation of LCE is that of Conti *et al.* [14], who did 3D finite element simulation based on the BTW model. They eliminated the orientation field  $\mathbf{n}$  from the BTW model by taking it to be the minimiser of the energy for fixed displacement field. The resulting energy as a functional of displacement field was non-convex, and so, they took its polyconvex envelope as the energy functional to analyse. Their simulation successfully recovered the stripe domain phenomenon. However, their model inherited the features of the BTW model, and only recovered ideally soft elasticity. In later work, the same authors applied a similar approach to the Verwey–Warner–Terentjev model and did successfully recover the semi-soft elasticity property [13].

In this paper, we extend the BTW energy by adding the Oseen–Frank energy and do finite element simulation for the full model on a 2D rectangular domain. We use the clamped pulling as a benchmark problem to check whether our numerical method can recover the special behaviours of LCE, such as stripe domain phenomenon and semi-soft elasticity.

The paper is organised as follows. In Section 2, we list the notations that we employ. In Section 3, we investigate the continuous problem and proceed to study the discrete one in Section 4. In Section 5, we present the numerical results, and in Section 6, we give the conclusions of the work. Finally in Section 7, we discuss possible improvements of current work to fully capture the stripe domain phenomenon.

## 2 Notations

In this paper, in addition to the standard notations of Sobolev spaces [1], such as  $W^{m,p}$ ,  $L^2$ ,  $H^1$ ,  $H_0^1$  and  $H^{-1}$ , we let

$$H_{0\Gamma}^1(\Omega) = \{u \in H^1(\Omega) : u = 0 \text{ on } \Gamma \subset \partial\Omega\}. \quad (2.1)$$

We use  $H_{g|\Gamma}^1(\Omega)$  to denote  $g + H_{0|\Gamma}^1(\Omega)$ , and  $\mathbf{H}_{g|\Gamma}^1(\Omega)$  to represent its vector version. We use  $H_{\Gamma}^{-1}(\Omega)$  to denote the dual space of  $H_{0|\Gamma}^1(\Omega)$ . We let  $\langle \cdot, \cdot \rangle$  represent the dot product of two vectors, the inner product in Hilbert spaces or the action of a linear functional on a function. Its actual meaning is made clear by the context.

We let  $\mathbb{M}^{m \times n}$  denote the space of real  $m \times n$  matrices. For any matrix  $F$ , we denote its transpose by  $F^T$ . For any square matrix  $A$ ,  $\det(A)$  denotes the determinant,  $\text{tr}(A)$  denotes the trace and  $\text{cof}(A)$  represents the cofactor matrix, whose  $(i, j)$  entry is equal to  $(-1)^{i+j}$  times the determinant of the submatrix obtained by eliminating row  $i$  and column  $j$  of the matrix  $A$ . For any matrix  $F$ , we use  $\frac{\partial \det}{\partial F}$  to denote  $\frac{\partial \det(F)}{\partial F}$ , which can be shown to be exactly  $\text{cof}(F)$ . We let  $A : B$  represent the inner product of the two matrices, that is

$$A : B = \text{tr}(A^T B) = \sum_{i,j} A_{ij} B_{ij}.$$

For any matrix  $F$ , we let  $|F|$  denote its Frobenius norm, that is

$$|F| = (F : F)^{1/2}.$$

### 3 Existence and well-posedness of the continuous problem

In this section, we present our analytical results on the continuous problem. We first introduce the energy functional and prove existence of minimiser. Then we derive the Euler–Lagrange equations for the energy and obtain the corresponding linearised system. Finally we reduce the linearised system to a standard saddle point framework and discuss its well-posedness.

#### 3.1 The energy functional and the minimisation problem

In this sub-section, we introduce the energy density as a combination of 2D BTW energy and the Oseen–Frank energy and prove existence of minimiser.

Throughout this paper, analysis and computation are carried out in 2D domains. This is justified as follows. In the experiment by Finkelmann *et al.* [25,35], the elastomer has a very thin, rectangular shape, and consequently, it can be assumed that the director vectors lie in the same plane. If the elastomer were compressed, it might buckle, with the directors tilting out of the plane. However, in the pulling experiment, the elastomer sample remains planar. If in addition, we look for configurations with director field confined in that plane, then a 2D model and analysis may be appropriate.

We use  $\mathbf{X} = (X_1, X_2)^T$  to denote a point in the reference configuration, and  $\mathbf{x}(\mathbf{X}) = (x_1, x_2)^T$  the corresponding point in the deformed configuration. We define the displacement field as  $\mathbf{u}(\mathbf{X}) = \mathbf{x} - \mathbf{X}$ . The deformation gradient tensor is  $F = \frac{\partial \mathbf{x}}{\partial \mathbf{X}}$  and satisfies  $F = I + \nabla \mathbf{u}$ . We assume the elastomer is incompressible, thus  $F$  satisfies  $\det(F) = 1$ . We denote the director field as  $\mathbf{n}(\mathbf{X}) = (n_1, n_2)^T$ , which is a unit vector representing the average orientation of the relevant liquid crystal molecular units at each point.

The BTW stored energy for LCE, in the form derived by DeSimone *et al.* [16–18], can be written as

$$W_{BTW} = \mu(|F|^2 - (1 - a)|F^T \mathbf{n}|^2 - 2a^{1/2}), \tag{3.1}$$

where  $\mu$  is an elasticity constant. The dimensionless constant  $a$  satisfies  $0 < a < 1$  and is a measurement of interaction between the bulk displacement  $\mathbf{u}$  and the director orientation  $\mathbf{n}$ . In the limit  $a \rightarrow 1$ , the BTW model degenerates to the neo-Hookean model, and there is no interaction between  $\mathbf{u}$  and  $\mathbf{n}$ . On the other hand, in the limit  $a \rightarrow 0$ , there is maximum interaction between  $\mathbf{u}$  and  $\mathbf{n}$ . Note that the reference configuration (the one with  $\mathbf{u} = 0$ ) for (3.1) is not the stress-free state. If we take the stress-free state as the reference state, the BTW energy will be in a slightly more complicated form. We will elaborate on this issue in the coming sections.

The Oseen–Frank stored energy [21], in its simplest form, can be written as

$$W_{OF} = b|\nabla\mathbf{n}|^2. \quad (3.2)$$

The energy (3.2) penalises change in the director field, and the prescribed constant  $b > 0$  measures the strength of the penalisation.

The non-dimensionalised energy functional is the following:

$$\begin{aligned} \Pi(\mathbf{u}, \mathbf{n}) = & \int_{\Omega} (|F|^2 - (1-a)|F^T\mathbf{n}|^2) + b|\nabla\mathbf{n}|^2 \\ & - \int_{\Omega} \mathbf{f} \cdot \mathbf{u} - \int_{\Gamma} \mathbf{g} \cdot \mathbf{u} da, \end{aligned} \quad (3.3)$$

where  $\mathbf{f}$  is a prescribed body force, and  $\mathbf{g}$  is an applied boundary traction on  $\Gamma \subset \partial\Omega$ . The admissible set for the displacement  $\mathbf{u}$  is

$$\mathcal{H} = \{\mathbf{u} \in H^1(\Omega, \mathbb{R}^2) : \det(I + \nabla\mathbf{u}) = 1 \text{ a.e. in } \Omega, \mathbf{u} = \mathbf{u}_0 \text{ on } \Gamma_u \subset \partial\Omega\}, \quad (3.4)$$

and, the admissible set for the director  $\mathbf{n}$  is

$$\mathcal{N} = \{\mathbf{n} \in H^1(\Omega, \mathbb{R}^2) : |\mathbf{n}| = 1 \text{ a.e. in } \Omega, \mathbf{n} = \mathbf{n}_0 \text{ on } \Gamma_n \subset \partial\Omega\}. \quad (3.5)$$

We let  $\mathcal{A} = \mathcal{H} \times \mathcal{N}$ . The admissible set  $\mathcal{A}$  is non-empty as long as  $\mathbf{u}_0$  and  $\mathbf{n}_0$  are both Lipschitz continuous functions [22].

The problem of energy minimisation is formulated as follows:

$$\text{Find } (\mathbf{u}, \mathbf{n}) \in \mathcal{A} \text{ minimising } \Pi(\mathbf{u}, \mathbf{n}) \text{ in } \mathcal{A}. \quad (3.6)$$

Next we prove existence of minimiser for (3.6).

Prior to the proof of existence, we summarise several lemmas, some of them well-known in the literature. We include them here for the purpose of self-completeness.

**Lemma 1** *Assume  $0 < a < 1$ ,  $|\mathbf{n}| = 1$  and  $\det(F) = 1$ . Then the BTW energy (3.1) is always non-negative. It is zero if and only if  $\text{eig}(FF^T) = \{a^{1/2}, a^{-1/2}\}$  and  $\mathbf{n}$  is an eigenvector corresponding to the eigenvalue  $a^{-1/2}$ .*

**Proof** Let the eigenvalues of  $F^T F$  be  $\lambda_1^2$  and  $\lambda_2^2$  and satisfy  $0 \leq \lambda_1^2 \leq \lambda_2^2$ , and let  $\mathbf{v}_1, \mathbf{v}_2$  denote the corresponding (unit) eigenvectors. Also assume that

$$\mathbf{n} = \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2.$$

Since  $|\mathbf{n}| = 1$ , and  $\mathbf{v}_1, \mathbf{v}_2$  are orthonormal, we have

$$\alpha_1^2 + \alpha_2^2 = 1.$$

Thus, we can rewrite the BTW energy as

$$\begin{aligned} W_{BTW} &= \mu(\text{tr}(FF^T) - (1 - a)\mathbf{n}^T FF^T \mathbf{n} - 2a^{1/2}) \\ &= \mu(\lambda_1^2 + \lambda_2^2 - (1 - a)(\alpha_1^2 \lambda_1^2 + \alpha_2^2 \lambda_2^2) - 2a^{1/2}). \end{aligned}$$

Since  $0 \leq \lambda_1^2 \leq \lambda_2^2$ , the values  $\alpha_2^2 = 1$  and  $\alpha_1^2 = 0$  minimise  $W_{BTW}$ . So,  $\mathbf{n}$  is parallel to  $\mathbf{v}_2$ . Consequently,

$$\begin{aligned} W_{BTW} &\geq \mu(\lambda_1^2 + \lambda_2^2 - (1 - a)\lambda_2^2 - 2a^{1/2}) \\ &= \mu(\lambda_1^2 + a\lambda_2^2 - 2a^{1/2}). \end{aligned}$$

Since  $\det(FF^T) = 1$ , we have that

$$\lambda_1^2 \lambda_2^2 = 1. \tag{3.7}$$

Therefore

$$\begin{aligned} W_{BTW} &\geq \mu(2\sqrt{\lambda_1^2 \cdot a\lambda_2^2} - 2a^{1/2}) \\ &= \mu(2a^{1/2} - 2a^{1/2}) \\ &= 0, \end{aligned}$$

hold. The equality is satisfied if and only if  $\lambda_1^2 = a\lambda_2^2$ . Combining it with (3.7) yields

$$\text{eig}(FF^T) = \{a^{1/2}, a^{-1/2}\}. \tag{3.8}$$

□

**Lemma 2** Assume  $|\mathbf{n}| = 1$  and  $0 < a < 1$ . Then

$$|F|^2 - (1 - a)|F^T \mathbf{n}|^2 \geq a|F|^2. \tag{3.9}$$

**Proof** Let  $\lambda_1^2$  and  $\lambda_2^2$  be as in Lemma 1. Since  $|\mathbf{n}| = 1$ , by the proof of Lemma 1, we have

$$\begin{aligned} |F|^2 - (1 - a)|F\mathbf{n}|^2 &\geq \lambda_1^2 + a\lambda_2^2 \\ &\geq a(\lambda_1^2 + \lambda_2^2) \\ &= a\text{tr}(FF^T) \\ &= a|F|^2. \end{aligned}$$

□

**Lemma 3** Assume  $|\mathbf{n}| = 1$ , and  $0 < a < 1$ . The function

$$L(F) = |F|^2 - (1 - a)|F^T \mathbf{n}|^2 \tag{3.10}$$

is convex with respect to  $F$ .

**Proof** Let

$$A(\mathbf{n}) = I - (1 - a)\mathbf{n}\mathbf{n}^T. \quad (3.11)$$

So,

$$L = \text{tr}(FAF^T). \quad (3.12)$$

For any matrices  $F_1, F_2 \in \mathbb{M}^{2 \times 2}$  and  $0 \leq \alpha \leq 1$ , we have

$$\begin{aligned} & [\alpha L(F_1) + (1 - \alpha)L(F_2)] - L(\alpha F_1 + (1 - \alpha)F_2) \\ &= \alpha \text{tr}(F_1 A F_1^T) + (1 - \alpha) \text{tr}(F_2 A F_2^T) \\ &\quad - \text{tr}[(\alpha F_1 + (1 - \alpha)F_2) A (\alpha F_1 + (1 - \alpha)F_2)^T] \\ &= \alpha(1 - \alpha) \text{tr}[(F_1 - F_2) A (F_1 - F_2)] \\ &= \alpha(1 - \alpha) (|F_1 - F_2|^2 - (1 - a)|(F_1 - F_2)^T \mathbf{n}|^2) \\ &\geq \alpha(1 - \alpha)a|F_1 - F_2|^2 \\ &\geq 0, \end{aligned}$$

where we have used Lemma 2. Hence, the result follows.  $\square$

Next we quote the following theorem on non-linear elasticity.

**Theorem 4** (Ball, [2]) *Let  $\Omega$  be a non-empty, bounded, open subset of  $\mathbb{R}^d$ .*

- *If  $d = 2$ , suppose we have  $\mathbf{u}_k \rightarrow \mathbf{u}$  in  $W^{1,s}$  with  $s > \frac{4}{3}$ , then we have  $\det(I + \nabla \mathbf{u}_k) \rightarrow \det(I + \nabla \mathbf{u})$  in  $\mathcal{D}'(\Omega)$ ;*
- *If  $d = 3$ ,*
  - *suppose we have  $\mathbf{u}_k \rightarrow \mathbf{u}$  in  $W^{1,s}$  with  $s > \frac{3}{2}$ , then we have  $\text{adj}(I + \nabla \mathbf{u}_k)_{ij} \rightarrow \text{adj}(I + \nabla \mathbf{u})_{ij}$  in  $\mathcal{D}'(\Omega)$ ;*
  - *suppose we have  $\mathbf{u}_k \rightarrow \mathbf{u}$  in  $W^{1,s}$ , and  $\text{adj}(I + \nabla \mathbf{u}_k) \rightarrow \text{adj}(I + \nabla \mathbf{u})$  in  $L^q(\Omega; \mathbb{M}^3)$  with  $s > 1, q > 1$  and  $\frac{1}{s} + \frac{1}{q} < \frac{4}{3}$ , then we have  $\det(I + \nabla \mathbf{u}_k) \rightarrow \det(I + \nabla \mathbf{u})$  in  $\mathcal{D}'(\Omega)$ .*

Now we are ready to prove the main theorem of this section.

**Theorem 5** *There exists solution to the problem (3.6).*

**Proof** Let  $m$  be the infimum of  $\Pi$  in  $\mathcal{A}$ , and let  $(\mathbf{u}_k, \mathbf{n}_k) \in \mathcal{A}$  be a minimising sequence of  $\Pi$ . Note that  $m < +\infty$ . Thus  $\Pi(\mathbf{u}_k, \mathbf{n}_k)$  is bounded above by some constant  $C$ . By

Lemma 2,

$$\begin{aligned}
 C \geq \Pi(\mathbf{u}_k, \mathbf{n}_k) &\geq \int_{\Omega} a|(I + \nabla \mathbf{u}_k)|^2 + b|\nabla \mathbf{n}_k|^2 dx \\
 &\quad - \|\mathbf{f}\|_{L^2(\Omega)} \|\mathbf{u}_k\|_{L^2(\Omega)} - \|\mathbf{g}\|_{L^2(\Gamma)} \|\mathbf{u}_k\|_{L^2(\Gamma)} \\
 &\geq \int_{\Omega} a|(I + \nabla \mathbf{u}_k)|^2 + b|\nabla \mathbf{n}_k|^2 dx \\
 &\quad - \left( \frac{1}{\varepsilon} \|\mathbf{f}\|_{L^2(\Omega)}^2 + \varepsilon \|\mathbf{u}_k\|_{L^2(\Omega)}^2 \right) - \left( \frac{1}{\varepsilon} \|\mathbf{g}\|_{L^2(\Gamma)}^2 + \varepsilon \|\mathbf{u}_k\|_{L^2(\Gamma)}^2 \right) \\
 &\geq \int_{\Omega} C_1|(I + \nabla \mathbf{u}_k)|^2 + C_2|\nabla \mathbf{n}_k|^2 dx - C_3,
 \end{aligned} \tag{3.13}$$

where  $\varepsilon > 0$  is small and  $C_i > 0, i = 1, 2, 3$  are constants. In the last step, we have applied the generalised Poincaré inequality ([9], p. 281) and the Trace Theorem ([20], p. 258). By (3.13),  $F_k = I + \nabla \mathbf{u}_k$  and  $\nabla \mathbf{n}_k$  are bounded in  $L^2$ . Since  $\nabla(\mathbf{u}_k - \mathbf{u}_0)$  is bounded in  $L^2$ , by the Poincaré inequality,  $\mathbf{u}_k$  is bounded in  $H^1$ . On the other hand, since  $\nabla \mathbf{n}_k$  is bounded in  $L^2$  and  $|\mathbf{n}_k| = 1$  a.e. in  $\Omega$ ,  $\mathbf{n}_k$  is bounded in  $H^1$ . Now since  $H^1$  is a reflexive Banach space and  $\mathbf{u}_k$  and  $\mathbf{n}_k$  are bounded in  $H^1$ , we can find a subsequence of  $\mathbf{u}_k$  and a subsequence of  $\mathbf{n}_k$  such that they are weakly convergent in  $H^1$ . We still denote them as  $(\mathbf{u}_k, \mathbf{n}_k)$ , and assume  $\mathbf{u}_k \rightharpoonup \mathbf{u}, \mathbf{n}_k \rightharpoonup \mathbf{n}$ .

Since  $\mathbf{u}_k \rightharpoonup \mathbf{u}$  in  $H^1$ , by Theorem 4,  $\det(I + \nabla \mathbf{u}_k) \rightarrow \det(I + \nabla \mathbf{u})$  in  $\mathcal{D}'(\Omega)$ . Moreover, since  $\det(I + \nabla \mathbf{u}_k) = 1$  a.e., it follows that  $\det(I + \nabla \mathbf{u}) = 1$  a.e. in  $\Omega$  as well<sup>1</sup>. On the other hand, weak convergence in  $H^1$  implies strong convergence in  $L^2$ , thus we can find a subsequence of  $\mathbf{n}_k$  that converges point-wise almost everywhere. Therefore we have  $|\mathbf{n}| = 1$  a.e. in  $\Omega$ . Finally since  $\mathbf{u}_k - \mathbf{u}_0 \in H^1_{0\Gamma_u}$ , which is a closed linear sub-space of  $H^1$ , by the Mazur's Theorem, it is weakly closed. Therefore  $\mathbf{u} - \mathbf{u}_0$  is also in  $H^1_{0\Gamma_u}$ , implying  $\mathbf{u} = \mathbf{u}_0$  on  $\Gamma_u$ . Similarly,  $\mathbf{n} = \mathbf{n}_0$  on  $\Gamma_n$  holds. Therefore  $(\mathbf{u}, \mathbf{n}) \in \mathcal{A}$ .

By Lemma 3, the following function:

$$L(F, \mathbf{n}, P) = (|F|^2 - (1 - a)|F^T \mathbf{n}|^2) + b|P|^2$$

<sup>1</sup> This is because, by definition,

$$\langle \det(I + \nabla \mathbf{u}_k), \phi \rangle \rightarrow \langle \det(I + \nabla \mathbf{u}), \phi \rangle$$

in  $\mathbb{R}$ , for any  $\phi \in \mathcal{D}(\Omega)$ . Since  $\det(I + \nabla \mathbf{u}_k) = 1$  a.e. in  $\Omega$  for any  $k$ ,

$$\langle \det(I + \nabla \mathbf{u}) - 1, \phi \rangle = 0, \quad \forall \phi \in \mathcal{D}(\Omega)$$

holds. Thus  $\det(I + \nabla \mathbf{u}) = 1$  a.e. in  $\Omega$ .

<sup>2</sup> This is because the embedding  $I : W^{1,p} \rightarrow L^p$  is compact for  $1 \leq p \leq \infty$  ([20], p. 274), while for any compact operator  $A : V \rightarrow W$  with  $V$  and  $W$  Banach spaces,  $u_k \rightharpoonup u$  in  $V$  implies  $Au_k \rightarrow Au$  in  $W$  ([9], Theorem 7.1-5 on p. 348).

is a convex function of  $F$  and  $P$ . Therefore by Theorem 1 of Section 8.2 of [20],  $\Pi$  is weakly lower semi-continuous. Thus

$$\begin{aligned} \Pi(\mathbf{u}, \mathbf{n}) &\leq \liminf_{k \rightarrow \infty} \Pi(\mathbf{u}_k, \mathbf{n}_k) \\ &= m. \end{aligned}$$

Since  $m$  is the infimum of  $\Pi$  on  $\mathcal{A}$ , we conclude that

$$\Pi(\mathbf{u}, \mathbf{n}) = m. \tag{3.14}$$

That is  $(\mathbf{u}, \mathbf{n})$  is the minimiser of  $\Pi$  on  $\mathcal{A}$ . □

### 3.2 Equilibrium equation and stress-free state

In this section, we derive the weak form of the equilibrium equation satisfied by the energy minimiser. After discretisation, this equilibrium equation can be regarded as a non-linear equation satisfied by the degrees of freedom (DOF) of the energy minimiser, which can then be solved by a non-linear solver, such as Newton’s method. We also discuss, in this sub-section, the solution corresponding to the stress-free state.

A common way to convert a constrained minimisation problem to an unconstrained one is by the method of the Lagrange multipliers. In this way, the constraints are moved from the admissible set to the objective function. We gain the flexibility at the cost of solving a larger system. After introducing the Lagrange multipliers, the objective energy functional becomes

$$\begin{aligned} \mathcal{E}(\mathbf{u}, \mathbf{n}, p, \lambda) &= \int_{\Omega} (|F|^2 - (1 - a)|F^T \mathbf{n}|^2) + b|\nabla \mathbf{n}|^2 \\ &\quad - p(\det F - 1) + \lambda(|\mathbf{n}|^2 - 1) - \int_{\Omega} \mathbf{f} \cdot \mathbf{u} - \int_{\Gamma} \mathbf{g} \cdot \mathbf{u}, \end{aligned} \tag{3.15}$$

where  $p \in L^2(\Omega)$  is the Lagrange multiplier for the incompressibility constraint  $\det(I + \nabla \mathbf{u}) = 1$ , which can be interpreted as *pressure*, while  $\lambda \in L^2(\Omega)$  is the Lagrange multiplier for the unity constraint  $|\mathbf{n}| = 1$ . The admissible set for  $(\mathbf{u}, \mathbf{n}, p, \lambda)$  is

$$\mathcal{S} = \mathbf{H}_{u_0|\Gamma_u}^1(\Omega) \times \mathbf{H}_{n_0|\Gamma_n}^1(\Omega) \times L^2(\Omega) \times L^2(\Omega). \tag{3.16}$$

Next we use variational principle to derive the weak form of the equilibrium equation, also known as the Euler–Lagrange equation. Assume  $(\mathbf{u}, \mathbf{n}, p, \lambda) \in \mathcal{S}$  minimises the energy (3.15). Then for any test function  $\mathbf{v} \in \mathbf{H}_{0|\Gamma_u}^1(\Omega)$ , the 1D function  $\mathcal{E}(\varepsilon) = \mathcal{E}(\mathbf{u} + \varepsilon \mathbf{v}, \mathbf{n}, p, \lambda)$  has a minimum at  $\varepsilon = 0$ . Thus it follows that

$$0 = \left. \frac{d\mathcal{E}}{d\varepsilon} \right|_{\varepsilon=0},$$



which simplifies to the following equation:

$$0 = \int_{\Omega} 2(F : \nabla v - (1 - a)\langle F^T \mathbf{n}, \nabla v^T \mathbf{n} \rangle) - p \frac{\partial \det}{\partial F} : \nabla v - \int_{\Omega} \mathbf{f} \cdot \mathbf{v} - \int_{\Gamma} \mathbf{g} \cdot \mathbf{v} da.$$

Similarly by taking the variations  $\mathbf{n} \rightarrow \mathbf{n} + \varepsilon \mathbf{m}$ ,  $p \rightarrow p + \varepsilon q$  or  $\lambda \rightarrow \lambda + \varepsilon \mu$ , we obtain the following Euler–Lagrange equations:

$$0 = \int_{\Omega} 2(F : \nabla v - (1 - a)\langle F^T \mathbf{n}, \nabla v^T \mathbf{n} \rangle) - p \frac{\partial \det}{\partial F} : \nabla v - \int_{\Omega} \mathbf{f} \cdot \mathbf{v} - \int_{\Gamma} \mathbf{g} \cdot \mathbf{v} da, \tag{3.17}$$

$$0 = \int_{\Omega} -2(1 - a)\langle F^T \mathbf{n}, F^T \mathbf{m} \rangle + 2b \nabla \mathbf{m} : \nabla \mathbf{n} + 2\lambda \langle \mathbf{n}, \mathbf{m} \rangle, \tag{3.18}$$

$$0 = \int_{\Omega} -q(\det F - 1), \tag{3.19}$$

$$0 = \int_{\Omega} \mu(\langle \mathbf{n}, \mathbf{n} \rangle - 1), \tag{3.20}$$

where the solution  $(\mathbf{u}, \mathbf{n}, p, \lambda)$  belongs to  $\mathcal{S}$ , while the test function  $(\mathbf{v}, \mathbf{m}, q, \mu)$  is in the space  $\mathbf{H}^1_{0|\Gamma_u} \times \mathbf{H}^1_{0|\Gamma_n} \times L^2(\Omega) \times L^2(\Omega)$ .

The corresponding equations in the strong form are derived using integration by parts, as long as  $\mathbf{u}$  and  $\mathbf{n}$  are smooth enough. The resulting system of partial differential equations is

$$\operatorname{div} \sigma + \mathbf{f} = 0 \quad \text{in } \Omega, \tag{3.21}$$

$$b \operatorname{div}(\nabla \mathbf{n}) + (1 - a)\mathbf{n}^T F F^T - \lambda \mathbf{n}^T = 0 \quad \text{in } \Omega, \tag{3.22}$$

$$\det F - 1 = 0 \quad \text{in } \Omega, \tag{3.23}$$

$$|\mathbf{n}|^2 - 1 = 0 \quad \text{in } \Omega, \tag{3.24}$$

$$\sigma \mathbf{v} = \mathbf{g} \quad \text{on } \partial\Omega \setminus \Gamma_u, \tag{3.25}$$

$$\frac{\partial \mathbf{n}}{\partial \mathbf{v}} = 0 \quad \text{on } \partial\Omega \setminus \Gamma_n, \tag{3.26}$$

where  $\mathbf{v}$  denotes the unit normal vector on the boundary, and

$$\sigma = 2(I - (1 - a)\mathbf{nn}^T)F - p \frac{\partial \det}{\partial F} \tag{3.27}$$

is the Piola–Kirchhoff stress tensor.

Note that when  $\mathbf{f} = 0$  and  $\mathbf{g} = 0$ , the reference configuration  $\mathbf{u} \equiv 0$  is not a stress-free state. The reason is as follows. At the reference configuration,  $F = I$ , so

$$\sigma = (2 - p)I - 2(1 - a)\mathbf{nn}^T. \tag{3.28}$$

The matrix  $(2 - p)I$  has rank 0 or 2 according to whether  $p = 2$  or  $p \neq 2$ , while the matrix  $2(1 - a)\mathbf{nn}^T$  has rank 1 for  $0 < a < 1$ . Thus the Cauchy stress  $\sigma$  cannot be zero<sup>3</sup>.

However, the stress-free state can be achieved by a uniform stretch of the reference state. It is easy to check that the system has zero stress in the case that

$$F = \begin{pmatrix} a^{1/4} & 0 \\ 0 & a^{-1/4} \end{pmatrix}, \tag{3.29}$$

$$\mathbf{n} \equiv (0, 1)^T, \quad p = 2\sqrt{a} \text{ and } \lambda = (1 - a)/\sqrt{a}.$$

### 3.3 Linearised system and well-posedness

In this section, we derive the linearised system of the equilibrium equations and discuss its well-posedness. This system is closely related to the matrix derivative in Newton’s method. Also, the well-posedness of the linearised system is closely related to the stability of the numerical scheme.

To linearise the original system, we fix a solution  $(\mathbf{u}, \mathbf{n}, p, \lambda)$ , and letting  $(\mathbf{w}, \mathbf{l}, o, \gamma)$  be a small perturbation, we carry out the corresponding Taylor expansions in equations (3.17)–(3.20) about the given solution, retaining the linear terms. The resulting linearised system is

$$a_1(\mathbf{w}, \mathbf{v}) + a_2(\mathbf{l}, \mathbf{v}) + b_1(o, \mathbf{v}) = L_1(\mathbf{v}), \tag{3.30}$$

$$a_2(\mathbf{m}, \mathbf{w}) + a_3(\mathbf{l}, \mathbf{m}) + b_2(\gamma, \mathbf{m}) = L_2(\mathbf{m}), \tag{3.31}$$

$$b_1(q, \mathbf{w}) = L_3(q), \tag{3.32}$$

$$b_2(\mu, \mathbf{l}) = L_4(\mu), \tag{3.33}$$

where  $a_1, a_2, a_3, b_1$  and  $b_2$  denote bi-linear forms depending on the solution  $(\mathbf{u}, \mathbf{n}, p, \lambda)$ , and  $L_1, L_2, L_3$  and  $L_4$  are linear functionals of the test functions. The perturbation  $(\mathbf{w}, \mathbf{l}, o, \gamma)$  is supposed to satisfy the linearised system (3.30)–(3.33) for any test function  $(\mathbf{v}, \mathbf{m}, q, \mu)$ . Both the perturbation and the test function belong to the space  $\mathbf{H}_{0|\Gamma_u}^1 \times \mathbf{H}_{0|\Gamma_n}^1 \times L^2(\Omega) \times H_{\Gamma_n}^{-1}$ . The bi-linear forms are defined by the following equations:

$$a_1(\mathbf{w}, \mathbf{v}) = \int_{\Omega} 2 \nabla \mathbf{w} : \nabla \mathbf{v} - 2(1 - a) \langle \nabla \mathbf{w}^T \mathbf{n}, \nabla \mathbf{v}^T \mathbf{n} \rangle - p \left( \frac{\partial^2 \det}{\partial F^2} \nabla \mathbf{w} \right) : \nabla \mathbf{v}, \tag{3.34}$$

$$a_2(\mathbf{m}, \mathbf{v}) = \int_{\Omega} -2(1 - a) \langle F^T \mathbf{m}, \nabla \mathbf{v}^T \mathbf{n} \rangle - 2(1 - a) \langle F^T \mathbf{n}, \nabla \mathbf{v}^T \mathbf{m} \rangle, \tag{3.35}$$

$$a_3(\mathbf{l}, \mathbf{m}) = \int_{\Omega} -2(1 - a) \langle F^T \mathbf{l}, F^T \mathbf{m} \rangle + 2b \nabla \mathbf{m} : \nabla \mathbf{l} + 2\lambda \langle \mathbf{l}, \mathbf{m} \rangle, \tag{3.36}$$

$$b_1(q, \mathbf{w}) = \int_{\Omega} -q \frac{\partial \det}{\partial F} : \nabla \mathbf{w}, \tag{3.37}$$

<sup>3</sup> If  $a = 1$ , and  $p = 2$ , we indeed get zero stress. In this case,  $p = 2$  corresponds to the hydrostatic pressure of a neo-Hookean material.

$$b_2(\mu, \mathbf{l}) = \int_{\Omega} 2\mu \langle \mathbf{l}, \mathbf{n} \rangle. \tag{3.38}$$

The linearised system (3.30)–(3.33) can now be reduced to a standard saddle point system. In fact, adding (3.30)–(3.31) together, and (3.32)–(3.33) together as well, yields

$$a(\tilde{\mathbf{w}}, \tilde{\mathbf{v}}) + b(\tilde{o}, \tilde{\mathbf{v}}) = \tilde{L}_1(\tilde{\mathbf{v}}), \tag{3.39}$$

$$b(\tilde{q}, \tilde{\mathbf{w}}) = \tilde{L}_2(\tilde{q}), \tag{3.40}$$

where  $\tilde{\mathbf{w}} = (\mathbf{w}, \mathbf{l})$ ,  $\tilde{\mathbf{v}} = (\mathbf{v}, \mathbf{m})$ ,  $\tilde{o} = (o, \gamma)$  and  $\tilde{q} = (q, \mu)$ . Moreover

$$a(\tilde{\mathbf{w}}, \tilde{\mathbf{v}}) = a_1(\mathbf{w}, \mathbf{v}) + a_2(\mathbf{l}, \mathbf{v}) + a_2(\mathbf{m}, \mathbf{v}) + a_3(\mathbf{l}, \mathbf{m}), \tag{3.41}$$

$$b(\tilde{q}, \tilde{\mathbf{v}}) = b_1(q, \mathbf{v}) + b_2(\mu, \mathbf{m}), \tag{3.42}$$

$$\tilde{L}_1(\tilde{\mathbf{v}}) = L_1(\mathbf{v}) + L_2(\mathbf{m}), \tag{3.43}$$

$$\tilde{L}_2(\tilde{q}) = L_3(q) + L_4(\mu). \tag{3.44}$$

The system (3.39) and (3.40) exhibits a saddle point structure. Its well-posedness is shown in the theorem by Ladyzenskaya–Babuska–Brezzi that we describe next as presented in [32].

**Theorem 6** (Ladyzenskaya–Babuska–Brezzi) *Consider the following saddle point problem:*

$$a(u, v) + b(p, v) = L_V(v) \quad \forall v \in \mathbf{V}, u \in \mathbf{V}, \tag{3.45}$$

$$b(q, u) = L_P(q) \quad \forall q \in \mathbf{P}, p \in \mathbf{P}, \tag{3.46}$$

with  $\mathbf{V}$  and  $\mathbf{P}$  given Hilbert spaces,  $L_V$  and  $L_P$  belonging to  $\mathbf{V}'$  and  $\mathbf{P}'$ , respectively. Moreover,  $a$  and  $b$  are continuous bi-linear forms defined on  $\mathbf{V} \times \mathbf{V}$  and  $\mathbf{P} \times \mathbf{V}$ , respectively. Define the operators

$$\mathcal{B} : \mathbf{V} \rightarrow \mathbf{P}'$$

$$v \mapsto \mathcal{B}v \text{ such that } \langle \mathcal{B}v, q \rangle = b(q, v) \quad \forall q \in \mathbf{P}$$

$$\mathcal{A} : \text{Ker } \mathcal{B} \rightarrow (\text{Ker } \mathcal{B})'$$

$$w \mapsto \mathcal{A}w \text{ such that } \langle \mathcal{A}w, v \rangle = a(w, v) \quad \forall v \in \text{Ker } \mathcal{B}$$

Then the operator  $\mathcal{B}$  is onto if and only if the spaces  $\mathbf{V}$  and  $\mathbf{P}$  satisfy the following inf-sup condition:

$$\inf_{q \in \mathbf{P}, \|q\|=1} \sup_{v \in \mathbf{V}, \|v\|=1} b(q, v) \geq \beta > 0. \tag{3.47}$$

Moreover, the mixed problem is well-posed if and only if  $\mathcal{B}$  is onto and  $\mathcal{A}$  is invertible.

According to the Ladyzenskaya–Babuska–Brezzi theorem, the well-posedness of the saddle point system requires both  $\mathcal{B}$  being onto and  $\mathcal{A}$  being invertible. Moreover  $\mathcal{B}$  being onto is equivalent to the inf-sup condition (3.47) being satisfied. What is the corresponding equivalent condition that guarantees the invertibility of  $\mathcal{A}$ ? It turns out that  $\mathcal{A}$  being invertible is also equivalent to an inf-sup condition, namely the inf-sup condition of  $a(\cdot, \cdot)$

on the space  $\text{Ker}\mathcal{B}$ ,

$$\inf_{w \in \text{Ker}\mathcal{B}, \|w\|=1} \sup_{v \in \text{Ker}\mathcal{B}, \|v\|=1} a(w, v) \geq \alpha > 0. \tag{3.48}$$

This can be easily proved using the Ladyzenskaya–Babuska–Brezzi theorem, and the fact that a linear operator  $A$  on a Hilbert space  $H$  is invertible if and only if  $A$  is onto and  $\text{Ker}(A) = 0$  ([31], p. 104). Therefore the well-posedness of a standard saddle point system amounts to the verification of the two inf-sup conditions (3.47) and (3.48). In practice, the inf-sup condition (3.48) is often replaced by the following stronger yet easier to verify *ellipticity* condition:

$$\inf_{v \in \text{Ker}\mathcal{B}, \|v\|=1} a(v, v) \geq \alpha > 0. \tag{3.49}$$

However, in most situations the spaces  $\mathbb{P}$  and  $\mathbb{V}$  are different, and so, the inf-sup condition (3.47) cannot be replaced with a stronger ellipticity condition, and it may be very difficult to verify analytically.

For the LCE problem that we study, the bi-linear form  $b(\tilde{q}, \tilde{v}) = b_1(q, v) + b_2(\mu, m)$  is the sum of two decoupled bi-linear forms. We prove that the inf-sup condition for  $b(\tilde{q}, \tilde{v})$  is actually equivalent to the inf-sup conditions for both  $b_1$  and  $b_2$ .

**Theorem 7** *The inf-sup condition for  $b(\tilde{q}, \tilde{v}) = b_1(q, v) + b_2(\mu, m)$  is satisfied if and only if the corresponding inf-sup conditions for  $b_1(q, v)$  and  $b_2(\mu, m)$  hold.*

**Proof** Assume the bi-linear form  $b_1(q, v)$  is defined on  $\mathbb{P} \times \mathbb{V}$  and  $b_2(\mu, m)$  is defined on  $A \times \mathbb{M}$ , where  $\mathbb{P}, \mathbb{V}, A, \mathbb{M}$  are Hilbert spaces.

First assume the inf-sup condition for  $b(\tilde{q}, \tilde{v})$  is satisfied. Then it follows from Theorem 6 that the operator

$$\begin{aligned} \mathcal{B} : \mathbb{V} \times \mathbb{M} &\rightarrow \mathbb{P}' \times A' \\ (v, m) &\mapsto \mathcal{B}(v, m) \text{ such that } \langle \mathcal{B}(v, m), (q, \mu) \rangle = b_1(q, v) + b_2(\mu, m) \quad \forall (q, \mu) \in \mathbb{P} \times A \end{aligned}$$

is onto. Therefore the operators

$$\begin{aligned} \mathcal{B}_1 : \mathbb{V} &\rightarrow \mathbb{P}' \\ v &\mapsto \mathcal{B}_1 v \text{ such that } \langle \mathcal{B}_1 v, q \rangle = b_1(q, v) \quad \forall q \in \mathbb{P} \end{aligned}$$

and

$$\begin{aligned} \mathcal{B}_2 : \mathbb{M} &\rightarrow A' \\ m &\mapsto \mathcal{B}_2 m \text{ such that } \langle \mathcal{B}_2 m, \mu \rangle = b_2(\mu, m) \quad \forall \mu \in A \end{aligned}$$

are both onto. Hence, it follows from Theorem 6 that the inf-sup conditions for  $b_1(q, v)$  and  $b_2(\mu, m)$  are satisfied.

Conversely let us assume the inf-sup conditions  $b_1(q, v)$  and  $b_2(\mu, m)$  are both satisfied. Then it follows that  $\mathcal{B}_1$  and  $\mathcal{B}_2$  are both onto, and so is the operator  $\mathcal{B}$ . Therefore by Theorem 6,  $b(\tilde{q}, \tilde{v}) = b_1(q, v) + b_2(\mu, m)$  satisfies the inf-sup condition.  $\square$

Consequently to verify the inf-sup condition for  $b(\tilde{q}, \tilde{\mathbf{v}}) = b_1(q, \mathbf{v}) + b_2(\mu, \mathbf{m})$ , it is sufficient to verify the inf-sup conditions for  $b_1(q, \mathbf{v})$  and  $b_2(\mu, \mathbf{m})$  individually. We point out that the bi-linear form  $b_1(q, \mathbf{v})$  in (3.37) corresponds to that of the incompressible elasticity problem [32], while  $b_2(\mu, \mathbf{m})$  in (3.38) corresponds to that of the harmonic map problem [24].

We observe that the inf-sup condition for  $b_1(q, \mathbf{v})$  is at least satisfied at the strain-free and the stress-free states. In fact, at the strain-free state,  $F = I$ , it reduces to that of the Stokes problem (for the proof, see for example [29]):

$$\inf_{q \in L^2(\Omega)} \sup_{\mathbf{v} \in \mathbf{H}_0^1(\Omega)} \frac{\langle q, \operatorname{div}(\mathbf{v}) \rangle}{\|q\|_0 \|\mathbf{v}\|_1} \geq \beta_1 > 0. \tag{3.50}$$

On the other hand, since the stress-free state has constant  $F$  matrix, the inf-sup condition for  $b_1(q, \mathbf{v})$  can be verified by change of variables [29]. In the general case that  $\mathbf{u} \neq 0$  and  $F$  is not a constant, analytical verification of such a condition can be very challenging.

The inf-sup condition for  $b_2(\mu, \mathbf{m})$  holds provided  $\mathbf{n}$  is sufficiently smooth. This is established in the following theorem, whose proof is a slight modification of that in [24]. The details can be found in [29].

**Theorem 8** *Assume  $\mathbf{n} \in \mathbf{H}_{n_0, \Gamma_n}^1(\Omega) \cap W^{1, \infty}(\Omega)$ , then the inf-sup condition for  $b_2(\mu, \mathbf{m})$  holds. That is*

$$\inf_{\mu \in H_{\Gamma_n}^1(\Omega)} \sup_{\mathbf{m} \in \mathbf{H}_{0, \Gamma_n}^1(\Omega)} \frac{\langle 2\mathbf{n} \cdot \mathbf{m}, \mu \rangle}{\|\mathbf{m}\|_1 \|\mu\|_{-1}} \geq \beta_2 > 0. \tag{3.51}$$

Finally to establish the ellipticity condition for the bi-linear form  $a(\tilde{\mathbf{w}}, \tilde{\mathbf{v}})$  is, in general, very complicated due to the complication of the expressions of  $a_1(\cdot, \cdot)$ ,  $a_2(\cdot, \cdot)$  and  $a_3(\cdot, \cdot)$ . We found that it actually does *not* hold at the stress-free state (see [29]). However, this does not imply that the linearised system (3.30)–(3.33) is ill-posed, since as previously mentioned, ellipticity is a sufficient condition instead of a necessary condition.

Although in many situations, the rigorous proof of the inf-sup conditions or ellipticity conditions is not available, the numerical ‘verification’ may be straightforward. We added quotation marks because numerical verification is not a rigorous argument and therefore cannot replace the analytical proof. However it may provide some insights when the analytical proof is not available. We will elaborate on this in later sections.

#### 4 Existence and well-posedness of the discrete problem

In this section, we investigate the existence and well-posedness of the discrete problem. Unlike the usual approach of simply replacing continuous spaces by finite element spaces, following Hu *et al.* [24], we include an interpolation operator in the discrete formulation. This operator plays an important role in the proof of existence and well-posedness of the discrete problem.

In this section, we first prove existence of minimiser of the discrete problem. We then derive the Euler–Lagrange equations and the corresponding linearised system. We also explain how to numerically compute the constants in the inf-sup and ellipticity

conditions as a way of verifying the well-posedness of the linearised system. Next we prove the existence and uniqueness of the Lagrange multipliers as a consequence of the inf-sup conditions for the discrete problem being satisfied. Finally we discuss some implementation issues, such as how to deal with the interpolation operator in the software package FEniCS, and how to numerically assemble the  $H^{-1}$  norm.

### 4.1 The discrete problem and existence of minimiser

As indicated in previous sections, the problem of finding  $(\mathbf{u}, p)$  is very similar to the case of the incompressible elasticity, while finding  $(\mathbf{n}, \lambda)$  bears an analogy with the harmonic map problem. Thus we choose the  $\mathbf{P}_2 \times P_1$  finite element spaces for  $(\mathbf{u}, p)$ , as in incompressible elasticity, and  $\mathbf{P}_1 \times P_1$  for  $(\mathbf{n}, \lambda)$ , following the harmonic map problem.

We let  $V_h$  denote the space of continuous piecewise linear functions and,  $V_{h,g|\Gamma} = \{v \in V_h \cap H^1 : v = g \text{ on } \Gamma\}$ . The symbols  $\mathbf{V}_h$  and  $\mathbf{V}_{h,g|\Gamma}$  refer to the corresponding vector version. We use  $\pi_h$  as the nodal interpolation operators onto the spaces  $V_h$  and  $\mathbf{V}_h$ . We let  $W_h$  denote the space of continuous piecewise quadratic functions and,  $W_{h,g|\Gamma} = \{w \in W_h \cap H^1 : w = g \text{ on } \Gamma\}$ . The symbols  $\mathbf{W}_h$  and  $\mathbf{W}_{h,g|\Gamma}$  denote the corresponding vector version as well.

The energy functional is still defined as

$$\begin{aligned} \Pi(\mathbf{u}, \mathbf{n}) &= \int_{\Omega} (|F|^2 - (1 - a)|F^T \mathbf{n}|^2) + b|\nabla \mathbf{n}|^2 \\ &\quad - \int_{\Omega} \mathbf{f} \cdot \mathbf{u} - \int_{\Gamma} \mathbf{g} \cdot \mathbf{u} da. \end{aligned} \tag{4.1}$$

We define the admissible set

$$\mathcal{A}_h = \mathcal{K}_h \times \mathcal{N}_h, \tag{4.2}$$

where

$$\mathcal{K}_h = \left\{ \mathbf{u}_h \in \mathbf{W}_{h,0|\Gamma_u} + \mathbf{u}_{0h}, \int_{\Omega} q_h(\det(I + \nabla \mathbf{u}_h) - 1) dx = 0, \forall q_h \in V_h \right\}, \tag{4.3}$$

and

$$\mathcal{N}_h = \left\{ \mathbf{n}_h \in \mathbf{V}_{h,0|\Gamma_n} + \mathbf{n}_{0h}, \int_{\Omega} \mu_h \pi_h(|\mathbf{n}_h|^2 - 1) dx = 0, \forall \mu_h \in V_{h,0|\Gamma_n} \right\}. \tag{4.4}$$

Notice that any piecewise linear function  $\mathbf{n}_h$  belongs to  $\mathcal{N}_h$  if and only if the function  $\pi_h(|\mathbf{n}_h|^2 - 1) \in V_{h,0|\Gamma_n}$  is identically 0, which means  $|\mathbf{n}_h| = 1$  at all the mesh nodes.

Our discrete formulation of the minimisation problem is

$$\text{Find } (\mathbf{u}_h, \mathbf{n}_h) \in \mathcal{A}_h, \text{ minimising } \Pi \text{ in } \mathcal{A}_h. \tag{4.5}$$

Before proving existence of minimiser, we first establish the following lemma.

**Lemma 9** Assume  $\mathbf{n} \in N_h$  and  $0 < a < 1$ , then for any matrix  $F \in \mathbb{M}^{2 \times 2}$

$$|F|^2 - (1 - a)|F^T \mathbf{n}|^2 \geq a|F|^2 \tag{4.6}$$

holds.

**Proof** Take any point  $x \in \Omega$ , and suppose it is inside the triangle  $\triangle P_1P_2P_3$ . Since  $\mathbf{n} \in N_h$ , we have

$$\mathbf{n}(x) = \lambda_1\mathbf{n}(P_1) + \lambda_2\mathbf{n}(P_2) + \lambda_3\mathbf{n}(P_3),$$

where  $\lambda_i, i = 1, 2, 3$  are barycentric coordinates. As indicated above,  $\mathbf{n} \in N_h$  if and only if  $|\mathbf{n}| = 1$  at all the mesh nodes. Thus it follows that

$$\begin{aligned} |\mathbf{n}(x)| &= |\lambda_1\mathbf{n}(P_1) + \lambda_2\mathbf{n}(P_2) + \lambda_3\mathbf{n}(P_3)| \\ &\leq \lambda_1|\mathbf{n}(P_1)| + \lambda_2|\mathbf{n}(P_2)| + \lambda_3|\mathbf{n}(P_3)| \\ &= \lambda_1 + \lambda_2 + \lambda_3 \\ &= 1. \end{aligned}$$

If  $|\mathbf{n}(x)| = 0$ , then the conclusion follows trivially. In the following, we assume that  $|\mathbf{n}(x)| > 0$ .

Let  $\hat{\mathbf{n}} = \mathbf{n}(x)/|\mathbf{n}(x)|$ , then  $|\hat{\mathbf{n}}| = 1$ . So

$$\begin{aligned} |F|^2 - (1 - a)|F^T\mathbf{n}|^2 &= |F|^2 - (1 - a)|\mathbf{n}(x)|^2|F^T\hat{\mathbf{n}}|^2 \\ &\geq |F|^2 - (1 - a)|F^T\hat{\mathbf{n}}|^2 \\ &\geq a|F|^2, \end{aligned}$$

where we have used Lemma 2 in the last step. □

Now we establish the following existence theorem.

**Theorem 10** *There exists a solution to the discrete minimisation problem (4.5).*

**Proof** Take any  $(\mathbf{u}_h, \mathbf{n}_h) \in \mathcal{A}_h$ . It follows from Lemma 9 that

$$\begin{aligned} \Pi(\mathbf{u}_h, \mathbf{n}_h) &\geq \int_{\Omega} a|(I + \nabla\mathbf{u}_h)|^2 + b|\nabla\mathbf{n}_h|^2 dx \\ &\quad - \|\mathbf{f}\|_{L^2(\Omega)}\|\mathbf{u}_h\|_{L^2(\Omega)} - \|\mathbf{g}\|_{L^2(\Gamma)}\|\mathbf{u}_h\|_{L^2(\Gamma)} \\ &\geq \int_{\Omega} a|(I + \nabla\mathbf{u}_h)|^2 + b|\nabla\mathbf{n}_h|^2 dx \\ &\quad - \left(\frac{1}{\varepsilon}\|\mathbf{f}\|_{L^2(\Omega)}^2 + \varepsilon\|\mathbf{u}_h\|_{L^2(\Omega)}^2\right) - \left(\frac{1}{\varepsilon}\|\mathbf{g}\|_{L^2(\Gamma)}^2 + \varepsilon\|\mathbf{u}_h\|_{L^2(\Gamma)}^2\right) \\ &\geq \int_{\Omega} C_1|(I + \nabla\mathbf{u}_h)|^2 + C_2|\nabla\mathbf{n}_h|^2 dx - C_3, \end{aligned}$$

where  $\varepsilon > 0$  is small, and  $C_i > 0, i = 1, 2, 3$  are constants. In the last step, we have applied the generalised Poincaré inequality ([9], p. 281) and the Trace Theorem ([20], p. 258). Thus  $\Pi(\mathbf{u}_h, \mathbf{n}_h) \rightarrow \infty$  as  $\|\mathbf{u}_h\|_1$  or  $\|\mathbf{n}_h\|_1$  goes to  $\infty$ . Hence its minimum must be achieved in a bounded subset of  $\mathcal{A}_h$ .

On the other hand, the admissible set  $\mathcal{A}_h$  is closed. The reason is as follows. Let  $\varphi_j, j = 1, \dots, N$  be a basis of  $V_h$ , and  $\psi_j, j = 1, \dots, M$  be a basis of  $V_{h,0|\Gamma_n}$ , and define

$$g_j(\mathbf{u}_h, \mathbf{n}_h) = \begin{cases} \int_{\Omega} \varphi_j(\det(I + \nabla \mathbf{u}_h) - 1) dx & 1 \leq j \leq N, \\ \int_{\Omega} \psi_{j-N} \pi_h(|\mathbf{n}_h|^2 - 1) dx & N + 1 \leq j \leq N + M. \end{cases} \tag{4.7}$$

Then  $g_j$  is a continuous function on  $(\mathbf{W}_{h,0|\Gamma_u} + \mathbf{u}_{0h}) \times (\mathbf{V}_{h,0|\Gamma_n} + \mathbf{n}_{0h})$ . Therefore  $\mathcal{A}_h$  can be written as the intersection of reciprocal images of 0 by the continuous functions  $g_j$ , so it is a closed set.

Since  $\Pi(\mathbf{u}_h, \mathbf{n}_h)$  is a continuous function on a closed, bounded *finite-dimensional* set, the Weierstrass Theorem guarantees the existence of  $(\mathbf{u}_h, \mathbf{n}_h) \in \mathcal{A}_h$  minimising  $\Pi$  in  $\mathcal{A}_h$ .  $\square$

### 4.2 Equilibrium equations and linearised system

Similar to the continuous problem, we convert the constrained minimisation problem to an unconstrained one by including the Lagrange multipliers. After that we derive the equilibrium equations and their linearisation. These equations are analogous to those of the continuous problem, except for the presence of the interpolation operator  $\pi_h$ .

After including the Lagrange multipliers, the discrete energy functional is given by

$$\begin{aligned} \mathcal{E}(\mathbf{u}, \mathbf{n}, p, \lambda) = & \int_{\Omega} (|F|^2 - (1 - a)|F^T \mathbf{n}|^2) + b|\nabla \mathbf{n}|^2 \\ & - p(\det(F) - 1) + \lambda(\pi_h \langle \mathbf{n}, \mathbf{n} \rangle - 1) \\ & - \int_{\Omega} \mathbf{f} \cdot \mathbf{u} - \int_{\Gamma} \mathbf{g} \cdot \mathbf{u} da, \end{aligned} \tag{4.8}$$

where  $F = I + \nabla \mathbf{u}$ .

Taking the first variation of the functional (4.8), we obtain the following equilibrium equations (Euler–Lagrange equations):

$$\begin{aligned} 0 = & \int_{\Omega} 2(F : \nabla \mathbf{v} - (1 - a)\langle F^T \mathbf{n}, \nabla \mathbf{v}^T \mathbf{n} \rangle) - p \frac{\partial \det}{\partial F} : \nabla \mathbf{v} \\ & - \int_{\Omega} \mathbf{f} \cdot \mathbf{v} - \int_{\Gamma} \mathbf{g} \cdot \mathbf{v} da, \end{aligned} \tag{4.9}$$

$$0 = \int_{\Omega} -2(1 - a)\langle F^T \mathbf{n}, F^T \mathbf{m} \rangle + 2b \nabla \mathbf{m} : \nabla \mathbf{n} + 2\lambda \pi_h \langle \mathbf{n}, \mathbf{m} \rangle, \tag{4.10}$$

$$0 = \int_{\Omega} -q(\det F - 1), \tag{4.11}$$

$$0 = \int_{\Omega} \mu \pi_h (\langle \mathbf{n}, \mathbf{n} \rangle - 1), \tag{4.12}$$

where the solution  $(\mathbf{u}, \mathbf{n}, p, \lambda) \in \mathbf{W}_{h,\mathbf{u}_0|\Gamma_u} \times \mathbf{V}_{h,\mathbf{n}_0|\Gamma_n} \times V_h \times V_{h,\lambda_0|\Gamma_n}$ , and the test function  $(\mathbf{v}, \mathbf{m}, q, \mu) \in \mathbf{W}_{h,0|\Gamma_u} \times \mathbf{V}_{h,0|\Gamma_n} \times V_h \times V_{h,0|\Gamma_n}$ .



Linearisation around a solution  $(\mathbf{u}, \mathbf{n}, p, \lambda)$  yields the system

$$a_1(\mathbf{w}, \mathbf{v}) + a_2(\mathbf{l}, \mathbf{v}) + b_1(o, \mathbf{v}) = L_1(\mathbf{v}), \tag{4.13}$$

$$a_2(\mathbf{m}, \mathbf{w}) + a_3(\mathbf{l}, \mathbf{m}) + b_2(\gamma, \mathbf{m}) = L_2(\mathbf{m}), \tag{4.14}$$

$$b_1(q, \mathbf{w}) = L_3(q), \tag{4.15}$$

$$b_2(\mu, \mathbf{l}) = L_4(\mu), \tag{4.16}$$

where both the perturbation  $(\mathbf{w}, \mathbf{l}, o, \gamma)$  and the test function  $(\mathbf{v}, \mathbf{m}, q, \mu)$  belong to  $\mathbf{W}_{h,0|\Gamma_u} \times \mathbf{V}_{h,0|\Gamma_n} \times V_h \times V_{h,0|\Gamma_n}$ . Here the bi-linear forms  $a_1, a_2, a_3, b_1$  and  $b_2$  depend on the solution  $(\mathbf{u}, \mathbf{n}, p, \lambda)$ . Moreover  $a_1, a_2$  and  $b_1$  are as in the continuous case, while  $a_3$  and  $b_2$  are slightly different and are given by

$$a_3(\mathbf{l}, \mathbf{m}) = \int_{\Omega} -2(1 - a)\langle F^T \mathbf{l}, F^T \mathbf{m} \rangle + 2b \nabla \mathbf{m} : \nabla \mathbf{l} + 2\lambda \pi_h \langle \mathbf{l}, \mathbf{m} \rangle, \tag{4.17}$$

and

$$b_2(\mu, \mathbf{m}) = \int_{\Omega} 2\mu \pi_h \langle \mathbf{n}, \mathbf{m} \rangle. \tag{4.18}$$

### 4.3 Well-posedness of the linearised system

As in the continuous case, verifying the well-posedness of the linearised system (4.13)–(4.16) can be reduced to verifying the inf-sup conditions for  $b_1(q, \mathbf{v}), b_2(\mu, \mathbf{m})$  and  $a(\tilde{\mathbf{w}}, \tilde{\mathbf{v}}) = a_1(\mathbf{w}, \mathbf{v}) + a_2(\mathbf{l}, \mathbf{v}) + a_2(\mathbf{m}, \mathbf{v}) + a_3(\mathbf{l}, \mathbf{m})$ .

The inf-sup condition for  $b_1(q, \mathbf{v})$  is also satisfied at least at the strain-free and stress-free state. In fact, at the strain-free state,  $F = I$ , this condition can be formulated as

$$\inf_{q \in V_h} \sup_{\mathbf{v} \in \mathbf{W}_{h,0|\Gamma_u}} \frac{\langle q, \text{div}(\mathbf{v}) \rangle}{\|q\|_0 \|\mathbf{v}\|_1} \geq \beta_1 > 0, \tag{4.19}$$

which is exactly the inf-sup condition for the Stokes problem with the Taylor–Hood element  $\mathbf{P}_2 \times P_1$  (proof of this inf-sup condition can be found, for example in Proposition 6.1 of [5]). The verification of the condition  $b_1(q, \mathbf{v})$  at the stress-free state follows as in the continuous case.

The proof of the inf-sup condition for  $b_2(\mu, \mathbf{m})$  is similar to the one in [24]. The only difference is that, in our case, the test functions  $\mu, \mathbf{m}$  are zero only on part of the boundary. A slight modification of the proof in [24] gives us the following. (The detailed proof can be found in [29].)

**Theorem 11** Assume  $\mathbf{n} \in \mathbf{H}_{n_0|\Gamma_n}^1(\Omega) \cap \mathcal{W}^{1,\infty}(\Omega)$ , and  $\mathbf{n}_h \in \mathbf{V}_{h,n_0|\Gamma_n}$  satisfies  $|\mathbf{n}_h| \geq C > 0$  and  $\|\mathbf{n}_h - \pi_h \mathbf{n}\|_1 \leq \gamma / |\log(h)|^{1/2}$ . Then there is a positive constant  $\beta_2$ , independent of  $h$ , such that

$$\inf_{\mu \in V_{h,0|\Gamma_n}} \sup_{\mathbf{m} \in \mathbf{V}_{h,0|\Gamma_n}} \frac{\langle \pi_h[\mathbf{n}_h \cdot \mathbf{m}], \mu \rangle}{\|\mu\|_{-1} \|\mathbf{m}\|_1} \geq \beta_2. \tag{4.20}$$

Theorem 11 states that, if the true solution  $\mathbf{n}$  is smooth, the approximate solution  $\mathbf{n}_h$  is close to it, and its norm is bounded below, then the inf-sup condition for  $b_2(\mu, \mathbf{m})$  is always satisfied.

For the inf-sup condition for  $a(\tilde{w}, \tilde{v})$  and the inf-sup conditions for  $b_1$  and  $b_2$  in general cases, analytical verification may turn out to be very difficult. However, in the discrete case, the inf-sup values and the ellipticity constants can be computed numerically. We can compute the inf-sup values for a series of finer and finer meshes. If these inf-sup values are bounded below by a positive constant, we infer some evidence that the inf-sup condition might be satisfied for all meshes. This type of verification known as *inf-sup test* [8] provides a convenient way to get information when analytical results are not available. However the inf-sup test cannot replace the analytical proof, because we cannot apply the test on infinite number of meshes.

For a general inf-sup condition, the inf-sup value

$$\beta_h = \inf_{q \in \mathbb{P}_h, \|q\|=1} \left\{ \sup_{v \in \mathbb{V}_h, \|v\|=1} b(q, v) \right\} \tag{4.21}$$

turns out to be related to the smallest singular value of certain matrix. The following theorem summaries the results from [5].

**Theorem 12** *Let the matrices  $S, T, B$  be defined by the following equations:*

$$\|q_h\|^2 = \mathbf{q}^T \mathbf{S} \mathbf{q}, \tag{4.22}$$

$$\|v_h\|^2 = \mathbf{v}^T \mathbf{T} \mathbf{v}, \tag{4.23}$$

$$b(q_h, v_h) = \mathbf{q}^T \mathbf{B} \mathbf{v}, \tag{4.24}$$

where  $\mathbf{q}, \mathbf{v}$  are the DOF of  $q_h$  and  $v_h$ , respectively. Then the inf-sup value  $\beta_h$  in (4.21) is equal to the smallest singular value of the matrix  $S^{-\frac{1}{2}} B T^{-\frac{1}{2}}$ .

In our case, we also want to compute the inf-sup value or ellipticity constant for the bi-linear form  $a(\cdot, \cdot)$  on  $\text{Ker}(\mathcal{B}_h)$ , where  $\mathcal{B}_h : \mathbb{V}_h \rightarrow \mathbb{P}'_h$  is defined by  $b(q, v) = (q, \mathcal{B}_h v)$  for any  $q \in \mathbb{P}_h$  and  $v \in \mathbb{V}_h$ . That is we want to compute the inf-sup value  $\hat{\beta}_h$  in

$$\hat{\beta}_h = \inf_{u \in \text{Ker}(\mathcal{B}_h), \|u\|=1} \sup_{v \in \text{Ker}(\mathcal{B}_h), \|v\|=1} a(u, v) \tag{4.25}$$

and the ellipticity constant  $\hat{\alpha}_h$  in

$$\hat{\alpha}_h = \inf_{v \in \text{Ker}(\mathcal{B}_h), \|v\|=1} a(v, v). \tag{4.26}$$

We prove the following result, which is similar to Theorem 12.

**Theorem 13** *Let  $n$  and  $m$  be the dimensions of  $\mathbb{V}_h$  and  $\mathbb{P}_h$ , respectively. Let the matrices  $T, A, B$  be defined by the following equations:*

$$\|v_h\|^2 = \mathbf{v}^T \mathbf{T} \mathbf{v},$$

$$a(u_h, v_h) = \mathbf{u}^T \mathbf{A} \mathbf{v},$$

$$v_h \in \text{Ker}(\mathcal{B}_h) \Leftrightarrow \mathbf{B} \mathbf{v} = \mathbf{0},$$

where  $\mathbf{u}, \mathbf{v}$  are the DOF of  $u_h$  and  $v_h$ , respectively. Assume that  $B$  is full rank, and let the matrix  $Q$  be defined by the QR decomposition of  $(BT^{-1/2})^T$

$$(BT^{-1/2})^T = Q \begin{pmatrix} R \\ 0 \end{pmatrix}.$$

Then the inf-sup value  $\hat{\beta}_h$  in (4.25) and the ellipticity constant  $\hat{\alpha}_h$  in (4.26) are, respectively, equal to the smallest singular value and the smallest eigenvalue of the matrix  $A_1$ , where  $A_1$  is the lower right  $(n - m) \times (n - m)$  submatrix of the matrix  $Q^T T^{-1/2} A T^{-1/2} Q$ .

**Proof** First, let  $\mathbf{x} = T^{\frac{1}{2}}\mathbf{u}, \mathbf{y} = T^{\frac{1}{2}}\mathbf{v}$ . So

$$\hat{\beta}_h = \inf_{\mathbf{x} \in \text{Ker}(\tilde{B})} \sup_{\mathbf{y} \in \text{Ker}(\tilde{B})} \frac{\mathbf{x}^T \tilde{A} \mathbf{y}}{\sqrt{\mathbf{x}^T \mathbf{x}} \sqrt{\mathbf{y}^T \mathbf{y}}}, \tag{4.27}$$

where  $\tilde{B} = BT^{-1/2}$ , and  $\tilde{A} = T^{-1/2} A T^{-1/2}$ .

Since  $\tilde{B}$  is full rank, the matrix  $R \in \mathbb{M}^{m \times m}$  in the QR decomposition

$$\tilde{B}^T = Q \begin{pmatrix} R \\ 0 \end{pmatrix} \tag{4.28}$$

is non-singular. Let

$$Q^T \mathbf{x} = \begin{pmatrix} \mathbf{w}_x \\ \mathbf{z}_x \end{pmatrix}, \tag{4.29}$$

where  $\mathbf{w}_x \in \mathbb{R}^m$  and  $\mathbf{z}_x \in \mathbb{R}^{n-m}$ . Then it is easy to verify that

$$\mathbf{x} \in \text{Ker}(\tilde{B}) \Leftrightarrow \mathbf{w}_x = 0.$$

Thus there is no constraint on  $\mathbf{z}_x$ . Therefore

$$\inf_{\mathbf{x} \in \text{Ker}(\tilde{B})} \sup_{\mathbf{y} \in \text{Ker}(\tilde{B})} \frac{\mathbf{x}^T \tilde{A} \mathbf{y}}{\sqrt{\mathbf{x}^T \mathbf{x}} \sqrt{\mathbf{y}^T \mathbf{y}}} = \inf_{\mathbf{z}_x} \sup_{\mathbf{z}_y} \frac{\mathbf{z}_x^T A_1 \mathbf{z}_y}{\sqrt{\mathbf{z}_x^T \mathbf{z}_x} \sqrt{\mathbf{z}_y^T \mathbf{z}_y}}, \tag{4.30}$$

where  $A_1$  is the lower right  $(n - m) \times (n - m)$  corner of the matrix  $Q^T \tilde{A} Q$ . Thus by Theorem 12,  $\hat{\beta}_h$  is equal to the smallest singular value of the matrix  $A_1$ .

Similarly we can show that  $\hat{\alpha}_h$  is equal to the smallest eigenvalue of the matrix  $A_1$ .  $\square$

*Remark* The matrix  $B$  is full-rank if and only if the operator  $\mathcal{B}_h$  is onto, which is true if and only if the following inf-sup condition holds

$$\inf_{q \in \mathbb{P}_h, \|q\|=1} \left\{ \sup_{v \in \mathbb{V}_h, \|v\|=1} b(q, v) \right\} \geq \beta_h > 0. \tag{4.31}$$

**4.4 Existence and uniqueness of the Lagrange multipliers for the discrete system**

For the incompressible elasticity problem, Le Tallec [27] proved existence and uniqueness of the Lagrange multiplier  $p$  given that the inf-sup condition for  $b_1(q, \mathbf{v})$  is satisfied. In this sub-section, we use similar arguments to prove existence and uniqueness of  $p$  and  $\lambda$  given that the inf-sup conditions for  $b_1(q, \mathbf{v})$  and  $b_2(\mu, \mathbf{m})$  are both satisfied. The proof uses the following result of Clarke [10, 11].

**Theorem 14** *Let  $J$  denote a finite set of integers. We suppose that the following are given:  $E$  a Banach space,  $g_0, g_j(j \in J)$  locally Lipschitz functions from  $E$  to  $\mathbb{R}$ , and  $C$  a closed subset of  $E$ . We consider the following problem:*

$$\begin{aligned} & \text{Minimise } g_0(x) \\ & \text{subject to } x \in C, \quad g_j(x) = 0, \quad \forall j \in J. \end{aligned} \tag{4.32}$$

If  $\bar{x}$  is a local solution of (4.32), then there exist real numbers  $r_0, s_j$  not all zero, and a point  $\xi$  in the dual space  $E'$  of  $E$  such that

$$\xi \in r_0 \partial g_0(\bar{x}) + \sum_j s_j \partial g_j(\bar{x}), \quad -\xi \in N_C(\bar{x}), \tag{4.33}$$

where  $N_C(\bar{x})$  is the normal cone at  $C$  in  $\bar{x}$ , and  $\partial g_j$  is the generalised gradient of  $g_j(x)$ .

Next we use Theorem 14 to prove the existence and uniqueness of  $p$  and  $\lambda$ .

**Theorem 15** *Suppose  $(\mathbf{u}_h, \mathbf{n}_h) \in \mathcal{X}_h \times \mathcal{N}_h$ , and at  $(\mathbf{u}_h, \mathbf{n}_h)$ , the inf-sup conditions for  $b_1$  and  $b_2$  are both satisfied. Then there exist a unique  $p_h \in V_h$  and a unique  $\lambda_h \in V_{h, \lambda_0|_{\Gamma_n}}$  such that  $(\mathbf{u}_h, \mathbf{n}_h, p_h, \lambda_h)$  is a solution of the discrete equilibrium equations (4.9)–(4.12).*

**Proof** Let us denote

$$E = C = (\mathbf{W}_{h,0|\Gamma_u} + \mathbf{u}_{0h}) \times (\mathbf{V}_{h,0|\Gamma_n} + \mathbf{n}_{0h}) \tag{4.34}$$

$$g_0(x) = \Pi(\mathbf{v}_h, \mathbf{m}_h), \quad g_j(x) = g_j(\mathbf{v}_h, \mathbf{m}_h), \tag{4.35}$$

where the functions  $g_j$  were defined in (4.7). It is easy to see that

$$N_C(\bar{x}) = N_E(\bar{x}) = (0, 0). \tag{4.36}$$

Notice that

$$\partial \Pi(\mathbf{u}_h, \mathbf{n}_h) \subset \{Dg_0^1 + Dg_0^2\}, \tag{4.37}$$

where  $Dg_0^1$  and  $Dg_0^2$  are in  $[(\mathbf{W}_{h,0|\Gamma_u} + \mathbf{u}_{0h}) \times (\mathbf{V}_{h,0|\Gamma_n} + \mathbf{n}_{0h})]^*$ . We have

$$Dg_0^1(\mathbf{u}_h, \mathbf{n}_h) \cdot (\mathbf{v}_h, \mathbf{m}_h) = \begin{pmatrix} f_1(\mathbf{v}_h) \\ f_2(\mathbf{m}_h) \end{pmatrix},$$

where

$$f_1(\mathbf{v}) = \int_{\Omega} 2(F_h : \nabla \mathbf{v} - (1 - a)\langle F_h^T \mathbf{n}_h, \nabla \mathbf{v}^T \mathbf{n}_h \rangle),$$

and

$$f_2(\mathbf{m}) = \int_{\Omega} -2(1 - a)\langle F_h^T \mathbf{n}_h, F_h^T \mathbf{m} \rangle + 2b \nabla \mathbf{m} : \nabla \mathbf{n}_h.$$

Also

$$Dg_0^2(\mathbf{u}_h, \mathbf{n}_h) \cdot (\mathbf{v}_h, \mathbf{m}_h) = \begin{pmatrix} -\int_{\Omega} \mathbf{f} \cdot \mathbf{v}_h - \int_{\Gamma} \mathbf{g} \cdot \mathbf{v}_h da \\ 0 \end{pmatrix}.$$

We also observe that  $g_j(\mathbf{v}_h, \mathbf{m}_h)$  is continuously differentiable in  $(\mathbf{W}_{h,0\Gamma_u} + \mathbf{u}_{0h}) \times (\mathbf{V}_{h,0\Gamma_n} + \mathbf{n}_{0h})$  and that

$$\partial g_j = \{Dg_j\}, \tag{4.38}$$

$$Dg_j(\mathbf{u}_h, \mathbf{n}_h) \cdot (\mathbf{v}_h, \mathbf{m}_h) = \begin{pmatrix} \int_{\Omega} -\varphi_j \frac{\partial \det}{\partial F}(I + \nabla \mathbf{u}_h) : \nabla \mathbf{v}_h \\ 0 \end{pmatrix}$$

for  $1 \leq j \leq N$ , (4.39)

$$Dg_j(\mathbf{u}_h, \mathbf{n}_h) \cdot (\mathbf{v}_h, \mathbf{m}_h) = \begin{pmatrix} 0 \\ \int_{\Omega} \psi_{j-N} \pi_h(2\mathbf{n}_h \cdot \mathbf{m}_h) dx \end{pmatrix}$$

for  $N + 1 \leq j \leq N + M$ . (4.40)

Therefore applying Theorem 14, we have that there exists real numbers  $r_0, s_j$ , not all zero, such that

$$0 \in r_0 \partial \Pi(\mathbf{u}_h, \mathbf{n}_h) + \sum_{j=1}^{N+M} s_j \partial g_j(\mathbf{u}_h, \mathbf{n}_h). \tag{4.41}$$

Using (4.37) and (4.38), equation (4.41) becomes

$$r_0 \{Dg_0^1(\mathbf{u}_h, \mathbf{n}_h) + Dg_0^2(\mathbf{u}_h, \mathbf{n}_h)\} + \sum_{j=1}^{N+M} s_j Dg_j(\mathbf{u}_h, \mathbf{n}_h) = 0,$$

in  $[(\mathbf{W}_{h,0\Gamma_u} + \mathbf{u}_{0h}) \times (\mathbf{V}_{h,0\Gamma_n} + \mathbf{n}_{0h})]^*$ . (4.42)

Assume now  $r_0 = 0$ . By the linearity property, and using (4.39) and (4.40), we rewrite (4.42) as follows:

$$\int_{\Omega} \left( \sum_{j=1}^N s_j \varphi_j \right) \frac{\partial \det}{\partial F}(I + \nabla \mathbf{u}_h) : \nabla \mathbf{v}_h dx = 0, \quad \forall \mathbf{v}_h \in \mathbf{W}_{h,0\Gamma_u}, \tag{4.43}$$

$$\int_{\Omega} \left( \sum_{j=1}^M s_{N+j} \psi_j \right) \pi_h(2\mathbf{n}_h \cdot \mathbf{m}_h) dx = 0, \quad \forall \mathbf{m}_h \in \mathbf{V}_{h,0\Gamma_n}. \tag{4.44}$$

Since at least one  $s_j$  is non-zero, at least one of the equations (4.43) or (4.44) is in contradiction with the inf-sup conditions. Thus  $r_0$  cannot be zero. We can then divide (4.42) by  $r_0$  to get

$$Dg_0^1(\mathbf{u}_h, \mathbf{n}_h) \cdot (\mathbf{v}_h, \mathbf{m}_h) + 1/r_0 \sum_{j=1}^M s_j Dg_j(\mathbf{u}_h, \mathbf{n}_h) = -Dg_0^2(\mathbf{u}_h, \mathbf{n}_h) \cdot (\mathbf{v}_h, \mathbf{m}_h),$$

$$\forall (\mathbf{v}_h, \mathbf{m}_h) \in [(\mathbf{W}_{h,0|\Gamma_u} + \mathbf{u}_{0h}) \times (\mathbf{V}_{h,0|\Gamma_n} + \mathbf{n}_{0h})]. \tag{4.45}$$

That is

$$f_1(\mathbf{v}_h) - \int_{\Omega} p_h \frac{\partial \det}{\partial F}(I + \nabla \mathbf{u}_h) : \nabla \mathbf{v}_h dx = \int_{\Omega} \mathbf{f} \cdot \mathbf{v}_h + \int_{\Gamma} \mathbf{g} \cdot \mathbf{v}_h da, \tag{4.46}$$

$$f_2(\mathbf{m}_h) + \int_{\Omega} \lambda_h \pi_h(2\mathbf{n}_h \cdot \mathbf{m}_h) dx = 0, \tag{4.47}$$

where we have denoted

$$p_h = \left( \sum_{j=1}^N s_j \varphi_j \right) / r_0, \tag{4.48}$$

$$\lambda_h = \left( \sum_{j=1}^M s_{N+j} \psi_j \right) / r_0. \tag{4.49}$$

Equations (4.46) and (4.47) are precisely (4.9) and (4.10). Since  $(\mathbf{u}_h, \mathbf{n}_h) \in \mathcal{K}_h \times \mathcal{N}_h$ , we conclude that  $(\mathbf{u}_h, \mathbf{n}_h, p_h, \lambda_h)$  is a solution of (4.9)–(4.12).

Finally if there were two distinct values  $p_h$ , their difference would violate the inf-sup condition for  $b_1$ . Likewise if there existed two distinct values  $\lambda_h$ , their difference would also violate the inf-sup condition for  $b_2$ . So we have the uniqueness of both  $p_h$  and  $\lambda_h$ . □

### 4.5 Some implementation issues

In this sub-section, we discuss issues related to the implementation of our numerical scheme, such as how to solve the non-linear problem using the software package FEniCS, how to deal with the interpolation operator  $\pi_h$  in FEniCS and how to assemble the  $H^{-1}$  norm when we assess the rate of convergence.

FEniCS [28] is an open source finite element package. It is very convenient to solve variational problems such as

$$a(u, v) = L(v), \quad \forall v \in \mathbb{V} \tag{4.50}$$

using FEniCS. To use FEniCS with C++, we just need to specify the finite element space  $\mathbb{V}_h$ , the expressions for the bi-linear form  $a(u, v)$  and the linear form  $L(v)$  in a form file such as ‘Poisson.uff’, and compile the form file into a C++ header file ‘Poisson.h’. Then in the C++ source file ‘main.cpp’, we specify the boundary conditions and let FEniCS proceed with the work, which includes assembling the matrix  $(a(\phi_i, \phi_j))$  and the right-hand

side ( $L(\phi_j)$ ), and calling solvers for the linear system<sup>4</sup>. To use FEniCS with Python is even simpler. The form file and the boundary conditions can be specified in the same Python script, and we can call FEniCS interactively in the Python shell.

Our problem (4.9)–(4.12), however, is a non-linear variational problem, and we cannot directly apply the above procedure in FEniCS. We explain here how to solve our non-linear problem using FEniCS. Let  $N$  be the dimension of the space  $\mathbf{W}_{h,0|\Gamma_u} \times \mathbf{V}_{h,0|\Gamma_n} \times V_h \times V_{h,0|\Gamma_n}$  for the test function  $(\mathbf{v}, \mathbf{m}, q, \mu)$ . Equations (4.9)–(4.12) can be regarded as a system of  $N$  non-linear equations for the DOF of the solution  $(\mathbf{u}, \mathbf{n}, p, \lambda)$ . We can solve it using a non-linear solver such as Newton’s method. It turns out that each iteration of Newton’s method is equivalent to solving the following linear variational problem for the increment  $(\mathbf{w}, \mathbf{l}, o, \gamma)$ :

$$a((\mathbf{w}, \mathbf{l}, o, \gamma), (\mathbf{v}, \mathbf{m}, q, \mu)) = L(\mathbf{v}, \mathbf{m}, q, \mu), \tag{4.51}$$

where the bi-linear form is

$$\begin{aligned} a((\mathbf{w}, \mathbf{l}, o, \gamma), (\mathbf{v}, \mathbf{m}, q, \mu)) = & a_1(\mathbf{w}, \mathbf{v}) + a_2(\mathbf{l}, \mathbf{v}) + a_2(\mathbf{m}, \mathbf{w}) + a_3(\mathbf{l}, \mathbf{m}) \\ & + b_1(o, \mathbf{v}) + b_1(q, \mathbf{w}) + b_2(\gamma, \mathbf{m}) + b_2(\mu, \mathbf{l}), \end{aligned} \tag{4.52}$$

and the linear form is

$$L(\mathbf{v}, \mathbf{m}, q, \mu) = -(F_1(\mathbf{v}) + F_2(\mathbf{m}) + F_3(q) + F_4(\mu)). \tag{4.53}$$

Here  $F_1, F_2, F_3$  and  $F_4$  are the right-hand sides of (4.9)–(4.12), respectively. The above observation can be verified by computing the derivative matrix and the right-hand side of Newton’s method, and comparing them with the matrix and the right-hand side of the above linear variational problem.

Another complication is that FEniCS does not support the interpolation operator  $\pi_h$  in their form file, at least not for the version 11.02 that we have used. We overcame this issue in the following way: we first let FEniCS assemble the matrix and the right-hand side without the  $\pi_h$  terms, then we manually assembled those terms and updated the matrix and the right-hand side. It turns out that we do not have to do numerical integration ourselves, instead we can compute those  $\pi_h$  terms using the DOF of  $\mathbf{n}_h$  and  $\lambda_h$ , and the matrix  $S = (\langle \varphi_i, \phi_j \rangle)$ , where the  $\varphi_i$ ’s denote the basis functions for the finite element space  $V_h$  of piecewise linear functions. The details can be found in [29].

Finally to compute the order of convergence, we need to compute the  $H^{-1}$  norm for any function in  $V_{h,0|\Gamma_n}$ . In the rest of this sub-section, we explain how to assemble the  $H^{-1}$  norm.

We first relate the  $H_{\Gamma_n}^{-1}$  norm of a function in  $V_{h,0|\Gamma_n}$  to the  $H^1$  norm of some other function in  $V_{h,0|\Gamma_n}$ . For any function  $v_h$  in  $V_{h,0|\Gamma_n}$ , we can define a linear functional  $g$  on  $H_{0|\Gamma_n}^1$  by

$$g(w) = \langle v_h, w \rangle_{L^2}, \quad \forall w \in H_{0|\Gamma_n}^1.$$

The  $H_{\Gamma_n}^{-1}$  norm of  $v_h$  is the same as the norm of the functional  $g$ . By the Riesz Representation Theorem, we can find  $v \in H_{0|\Gamma_n}^1$  such that

$$g(w) = \langle w, v \rangle_{H^1}, \quad \forall w \in H_{0|\Gamma_n}^1.$$

<sup>4</sup> Here we have used  $\phi_i$  to denote basis function of the finite element space  $\mathbf{W}_h$ .

Thus the norm of  $g$  is just the  $H^1$  norm of  $v$ . Therefore we get

$$\|v_h\|_{H_{\Gamma_n}^{-1}} = \|v\|_{H^1}. \tag{4.54}$$

Let  $\hat{v}_h$  be the  $L^2$  projection of  $v$  into  $V_{h,0|\Gamma_n}$ , then the  $H^1$  norm of  $v$  can be approximated by the  $H^1$  norm of  $\hat{v}_h$ .

Next we explain how to calculate the  $H^1$  norm of  $\hat{v}_h$ , which can be used to approximate the  $H_{\Gamma_n}^{-1}$  norm of  $v_h$ . Let  $\{\varphi_i, i = 1, \dots, n\}$  be a basis of  $V_{h,0|\Gamma_n}$ . We want to assemble the matrix  $S$  such that

$$\|v_h\|_{H_{\Gamma_n}^{-1}} \approx \|\hat{v}_h\|_{H^1} = \mathbf{v}^T S \mathbf{v},$$

where  $\mathbf{v} \in \mathbb{R}^n$  is the DOF for  $v_h$ .

**Theorem 16** *Let  $A$  and  $B$  be the matrices that satisfy*

$$\begin{aligned} \|v_h\|_{L^2} &= \mathbf{v}^T A \mathbf{v}, \\ \|v_h\|_{H^1} &= \mathbf{v}^T B \mathbf{v}, \end{aligned}$$

for any  $v_h$  in  $V_{h,0|\Gamma_n}$ , where  $\mathbf{v} \in \mathbb{R}^n$  is the DOF for  $v_h$ . Then the matrix  $S = AB^{-1}A$ .

**Proof** Let  $f : H_{\Gamma_n}^{-1} \rightarrow V_{h,0|\Gamma_n}$  be the map taking any  $v_h \in V_{h,0|\Gamma_n}$  to  $\hat{v}_h \in V_{h,0|\Gamma_n}$ , and let  $\hat{\varphi}_i = f(\varphi_i)$ . It is easy to see that

$$S_{ij} = \langle \hat{\varphi}_i, \hat{\varphi}_j \rangle_{H^1}. \tag{4.55}$$

By definition of  $\hat{\varphi}_i$ , we have

$$\int \varphi_i \varphi_j = \int D \hat{\varphi}_i D \varphi_j + \int \hat{\varphi}_i \varphi_j \quad \forall 1 \leq i, j \leq n. \tag{4.56}$$

Since  $\hat{\varphi}_i \in V_{h,0|\Gamma_n}$ , we can write

$$\hat{\varphi}_i = \sum_k G_{ik} \varphi_k.$$

Substituting it into (4.56) gives

$$\int \varphi_i \varphi_j = \sum_k G_{ik} \left( \int D \varphi_k \cdot D \varphi_j + \int \varphi_k \varphi_j \right).$$



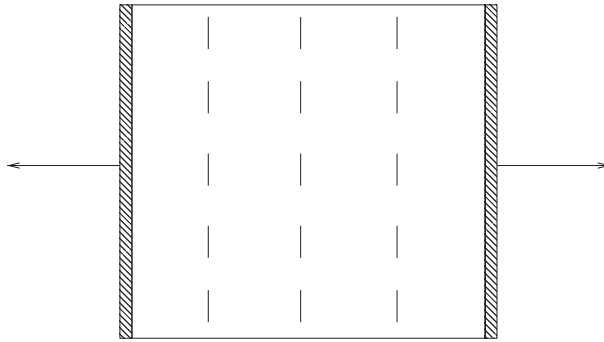


FIGURE 1. The elastomer is clamped and pulled on both sides.

That is  $A = GB$ , or  $G = AB^{-1}$ . Therefore

$$\begin{aligned}
 S &= (\langle \hat{\varphi}_i, \hat{\varphi}_j \rangle_{H^1}) \\
 &= (\langle D\hat{\varphi}_i, D\hat{\varphi}_j \rangle + \langle \hat{\varphi}_i, \hat{\varphi}_j \rangle) \\
 &= \left( \sum_{p,q} G_{ip} G_{jq} [\langle D\varphi_p, D\varphi_q \rangle + \langle \varphi_p, \varphi_q \rangle] \right) \\
 &= \left( \sum_{p,q} G_{ip} B_{pq} G_{qj}^T \right) \\
 &= GBG^T \\
 &= (AB^{-1})B(B^{-1}A) \\
 &= AB^{-1}A.
 \end{aligned}$$

□

*Remark* For any  $v_h$  in  $V_{h,0|\Gamma_n}$ , although  $v^T S v$  only gives approximate estimate of its  $H^{-1}$  norm, the difference goes to zero when  $h$  goes to 0.

### 5 Numerical results

In this section, we present results of the numerical simulation of the clamped-pulling experiment.

The simulation setup is as follows (Figure 1). The LCE is initially rectangular shaped and the directors align in the vertical direction. It is then clamped on the left and right edges and pulled in the horizontal direction.

As previously pointed out, in our model the stress-free state is different from the reference state. In the model and subsequent computation,  $\mathbf{u}$  represents the displacement relative to the reference domain. However the LCE should be in the stress-free state before it is clamped and pulled. This is not a big issue, because the stress-free state actually has constant deformation gradient matrix  $F$ , which means that it can be achieved by a uniform stretch from the reference state. Thus at both the reference and the stress-free

states, the LCE is rectangular, and displacements relative to the reference state or the stress-free state can be easily converted into each other. We take our reference domain to be the rectangle  $[0, L] \times [0, 1]$ . It can be verified that for the aspect ratio at the stress-free state to be AR, we should take

$$L = \frac{1}{\sqrt{a}}AR. \tag{5.1}$$

We give the following starting values for  $(\mathbf{u}, \mathbf{n}, p, \lambda)$  so that they correspond to the stress-free state:

$$u_X = (a^{1/4} - 1)(X - 0.5L), \tag{5.2}$$

$$u_Y = (a^{-1/4} - 1)(Y - 0.5), \tag{5.3}$$

$$\mathbf{n} \equiv (0, 1)^T, \tag{5.4}$$

$$p = 2\sqrt{a}, \tag{5.5}$$

$$\lambda = (1 - a)/\sqrt{a}, \tag{5.6}$$

where  $u_X$  and  $u_Y$  are the components of  $\mathbf{u}$ .

The physics of ‘clamped-pulling’ can be modelled by the following boundary conditions: at the two clamped edges,  $u_Y$  and  $\mathbf{n}$  remain at the starting values, while  $u_X$  decreases or increases uniformly (that is, independent of  $Y$ ). Although our model was not formulated as a time-dependent problem, we can still obtain information on the dynamical behaviour by solving a series of static problems, each of which only differ slightly from the previous one in the  $u_X$  boundary condition.

Notice that the problem is completely symmetric about the two centre lines  $X = 0.5L$  and  $Y = 0.5$ . Therefore we only need to do the computation on the upper-right quarter of the reference domain. The solution on the rest of the domain can be obtained by reflection.

Based on the discussion above, we list here the boundary conditions on the computation domain  $[0.5L, L] \times [0.5, 1]$ . First, to model the clamped-pulling set up, we impose the following Dirichlet boundary conditions at the clamped edge  $X = L$ :

$$u_X = 0.5L[a^{1/4}(1 + Mt) - 1], \tag{5.7}$$

$$u_Y = (a^{-1/4} - 1)(Y - 0.5), \tag{5.8}$$

$$\mathbf{n} = (0, 1)^T. \tag{5.9}$$

That is both  $u_Y$  and  $\mathbf{n}$  remain at their starting values, while  $u_X$  varies with  $t$ . Here  $t \in [0, 1]$  is the percentage of the loading. When  $t = 1$ , the LCE reaches its maximum elongation  $1 + M$ , where *elongation* is defined as the current length divided by the starting length (length at the stress-free state). Note that, by symmetry, the vertical centre line remains at  $X = 0.5L$ , while the horizontal centre line stays at  $Y = 0.5$ . Also by symmetry, the directors at these centre lines must be either strictly vertical or strictly horizontal. We assume that the directors change continuously during the pulling process, thus the directors at the two centre lines must stay at their starting values. Therefore we impose

Table 1. The numerical errors and orders of convergence

|                      | AR = 1   |          |          | AR = 3   |          |          |
|----------------------|----------|----------|----------|----------|----------|----------|
|                      | N = 8    | N = 16   | N = 32   | N = 8    | N = 16   | N = 32   |
| $\ e_u\ _0$          | 8.39E-04 | 2.69E-04 | 6.99E-05 | 1.51E-03 | 5.18E-04 | 1.41E-04 |
| order                | -        | 1.64     | 1.95     | -        | 1.55     | 1.88     |
| $\ e_u\ _1$          | 2.02E-02 | 7.66E-03 | 3.32E-03 | 2.14E-02 | 8.35E-03 | 3.57E-03 |
| order                | -        | 1.40     | 1.21     | -        | 1.36     | 1.23     |
| $\ e_n\ _0$          | 3.05E-02 | 8.25E-03 | 2.12E-03 | 3.79E-02 | 1.16E-02 | 3.20E-03 |
| order                | -        | 1.88     | 1.96     | -        | 1.71     | 1.86     |
| $\ e_n\ _1$          | 1.19E+00 | 6.23E-01 | 3.14E-01 | 1.48E+00 | 7.64E-01 | 3.81E-01 |
| order                | -        | 0.93     | 0.99     | -        | 0.95     | 1.00     |
| $\ e_p\ _0$          | 2.38E-02 | 8.99E-03 | 3.34E-03 | 2.14E-02 | 8.22E-03 | 2.79E-03 |
| order                | -        | 1.41     | 1.43     | -        | 1.38     | 1.56     |
| $\ e_\lambda\ _{-1}$ | 1.51E-03 | 5.22E-04 | 1.70E-04 | 1.95E-03 | 6.15E-04 | 1.92E-04 |
| order                | -        | 1.54     | 1.62     | -        | 1.66     | 1.68     |

the following boundary conditions at the two centre lines:

$$u_X = 0 \quad \text{on } X = 0.5L, \tag{5.10}$$

$$u_Y = 0 \quad \text{on } Y = 0.5, \tag{5.11}$$

$$\mathbf{n} = (0, 1)^T \quad \text{on } X = 0.5L \text{ and } Y = 0.5. \tag{5.12}$$

Finally to ensure that  $|\mathbf{n}| = 1$  at all the mesh nodes, we need to impose Dirichlet boundary condition for  $\lambda$  on the same boundary as  $\mathbf{n}$ . Thus the boundary condition for  $\lambda$  is

$$\lambda = (1 - a)/\sqrt{a} \quad \text{on } X = 0.5L, X = L \text{ and } Y = 0.5. \tag{5.13}$$

In our computation, we take  $a = 0.6$ ,  $b = 0.0015$  and  $M = 0.4$ . We slowly increase ‘load’  $t$  from 0 to 1 in a step size  $\Delta t = 0.01$ . We take the initial aspect ratio AR to be either 1 or 3. We use uniform mesh of size  $(AR \cdot N) \times N$ , where  $N$  is an integer. Each small rectangle of the mesh contains two triangles, which are split by the lower-left to upper-right diagonal.

Table 1 lists the numerical errors and orders of convergence. Here  $e_u$ ,  $e_n$ ,  $e_p$  and  $e_\lambda$  are the numerical errors for  $\mathbf{u}$ ,  $\mathbf{n}$ ,  $p$  and  $\lambda$ , respectively. And  $\|\cdot\|_0$ ,  $\|\cdot\|_1$  and  $\|\cdot\|_{-1}$  represent the  $L^2$ ,  $H^1$  and  $H^{-1}$  norms, respectively. The numerical errors are calculated for the solutions of adjacent meshes. For instance let us consider  $e_u$ . We first compute the solution  $\mathbf{u}^{(N)}$  and  $\mathbf{u}^{(2N)}$  on the mesh  $N$  and  $2N$ , respectively, next we interpolate the solution  $\mathbf{u}^{(N)}$  to the mesh  $2N$ , and finally, we compute the difference of that interpolation with the solution  $\mathbf{u}^{(2N)}$  and obtain  $e_u$ . The order of convergence is calculated in the usual sense. Take  $\|e_u\|_0$ , for example the order of convergence is calculated by

$$\frac{\log (\|e_u\|_0^{(N/2)} / \|e_u\|_0^{(N)})}{\log (2)}.$$

Table 2. The inf-sup values and ellipticity constants

| $t = 0$   | AR = 1    |           |           | AR = 3    |           |           |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
|           | $N = 4$   | $N = 8$   | $N = 16$  | $N = 4$   | $N = 8$   | $N = 16$  |
| $\beta_1$ | 0.5875    | 0.5879    | 0.5880    | 0.5877    | 0.5879    | 0.5880    |
| $\beta_2$ | 2.0000    | 2.0000    | 2.0000    | 2.0000    | 2.0000    | 2.0000    |
| $\beta$   | 2.70E-04  | 5.69E-05  | 1.62E-04  | 1.61E-04  | 1.21E-05  | 3.60E-05  |
| $\alpha$  | -1.27E-02 | -1.78E-02 | -1.21E-02 | -7.87E-02 | -5.42E-02 | -5.32E-02 |
| $t = 1$   | $N = 4$   | $N = 8$   | $N = 16$  | $N = 4$   | $N = 8$   | $N = 16$  |
| $\beta_1$ | 0.6431    | 0.6287    | 0.6163    | 0.6229    | 0.6125    | 0.6025    |
| $\beta_2$ | 1.9503    | 1.9065    | 1.8711    | 1.8737    | 1.7804    | 1.7517    |
| $\beta$   | 1.20E-03  | 5.82E-04  | 4.88E-05  | 3.48E-04  | 1.46E-04  | 2.55E-04  |
| $\alpha$  | -2.58E-03 | -5.82E-04 | -4.88E-05 | -2.50E-03 | -3.09E-03 | -3.34E-03 |

Table 1 shows that the  $L^2$  errors of  $\mathbf{u}$  and  $\mathbf{n}$  converge at rates close to 2, while their  $H^1$  errors converge at rates close to 1. The  $L^2$  error of  $p$  and the  $H^{-1}$  error of  $\lambda$  converge at rates of at least 1.

Table 2 lists the inf-sup values and the ellipticity constants at both the initial state ( $t = 0$ ) and at the final state ( $t = 1$ ) of the pulling process. Here  $\beta_1$  and  $\beta_2$  are the inf-sup values of the bi-linear forms  $b_1(\cdot, \cdot)$  and  $b_2(\cdot, \cdot)$ , respectively, while  $\beta$  and  $\alpha$  are the inf-sup value and ellipticity constant, respectively, of the bi-linear form  $a(\cdot, \cdot)$  on  $\text{Ker}(\mathcal{B})$ . The eigenvalue decomposition, singular value decomposition and QR decomposition were done using the open source library ALGLIB 2.6 [4]. Notice that for all cases in Table 2, the  $\alpha$ 's are negative, while  $\beta_1$ ,  $\beta_2$  and  $\beta$ 's are all positive. This means that, in all these cases, although the ellipticity conditions for  $a(\cdot, \cdot)$  are not satisfied, the inf-sup conditions for  $b_1(\cdot, \cdot)$ ,  $b_2(\cdot, \cdot)$  and  $a(\cdot, \cdot)$  are all satisfied, and therefore, the linearised system is well-posed. Furthermore, for both  $t = 0$  and  $t = 1$ , the inf-sup values  $\beta_1$  and  $\beta_2$  do not seem to change very much as the mesh refines. This suggests that the inf-sup values for  $b_1(\cdot, \cdot)$  and  $b_2(\cdot, \cdot)$  might have a constant positive lower bound during the whole pulling process, for all uniform meshes. On the other hand, this is not the case for the inf-sup value  $\beta$ . There is no obvious constant positive lower bound for  $\beta$ . In the case that  $\text{AR} = 1$  and  $t = 1$ , the  $\beta$  values even seem to go to zero as the mesh keeps on refining.

Next we check the stress-strain curve for semi-soft elasticity. Figures 2 and 3 show the stress-strain curves for  $\text{AR} = 1$  and  $\text{AR} = 3$ . In these figures, the  $x$ -axis is the strain, which is calculated by  $Mt$ , while the  $y$ -axis is the nominal stress, which is calculated by

$$\int_{\Gamma} \sigma(t) \mathbf{v} \cdot \mathbf{v} da, \tag{5.14}$$

where  $\Gamma$  is the clamped edge  $X = L$ , and  $\mathbf{v}$  is the normal vector on  $\Gamma$ . In both figures, the LCE is first hard, then soft, then hard again. Therefore we have successfully recovered the semi-soft elasticity.

To check how the soft regime changes with the meshes, we list in Table 3 the endpoints of the soft regime. Since the hard regimes have relatively small curvature, while the soft regimes have relatively large curvature, we choose the endpoints of the soft regime to be

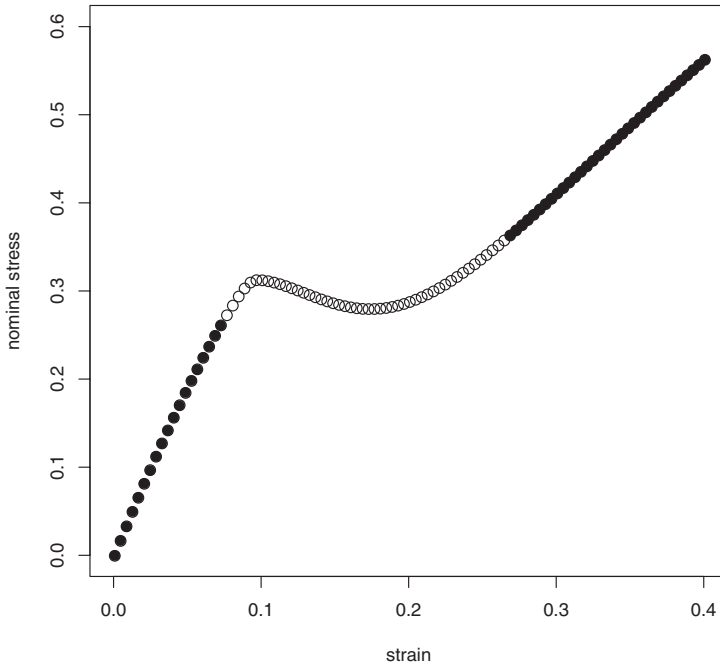


FIGURE 2. Nominal stress versus strain for  $AR = 1$ . Mesh size  $N = 32$ . The empty circles correspond to the soft regime  $[0.076, 0.264]$ , which is determined using a curvature criteria.

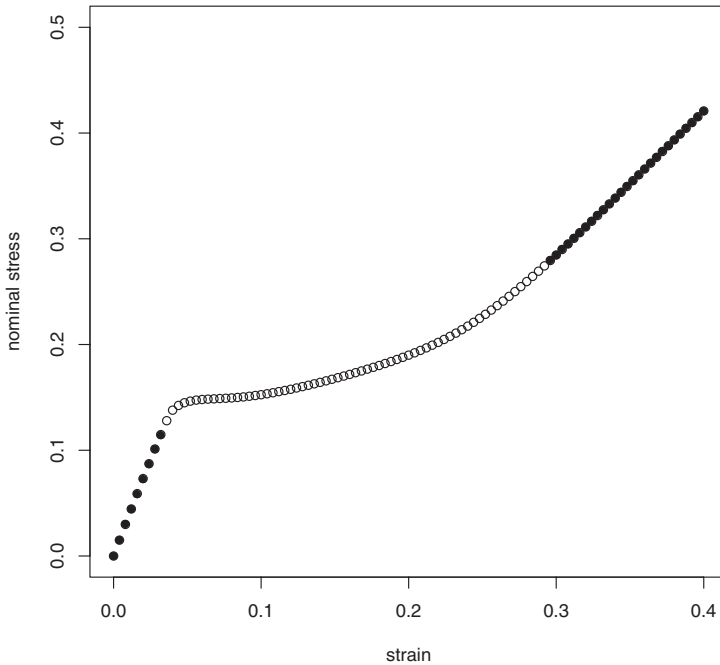


FIGURE 3. Nominal stress versus strain for  $AR = 3$ . Mesh size  $N = 32$ . The empty circles correspond to the soft regime  $[0.036, 0.292]$ , which is determined using a curvature criteria.

Table 3. *The endpoints of the soft regime*

|       | AR = 1 |       |        |        | AR = 3 |       |        |        |
|-------|--------|-------|--------|--------|--------|-------|--------|--------|
|       | N = 4  | N = 8 | N = 16 | N = 32 | N = 4  | N = 8 | N = 16 | N = 32 |
| left  | 0.096  | 0.076 | 0.076  | 0.076  | 0.048  | 0.040 | 0.036  | 0.036  |
| right | 0.288  | 0.272 | 0.264  | 0.264  | 0.276  | 0.288 | 0.292  | 0.292  |

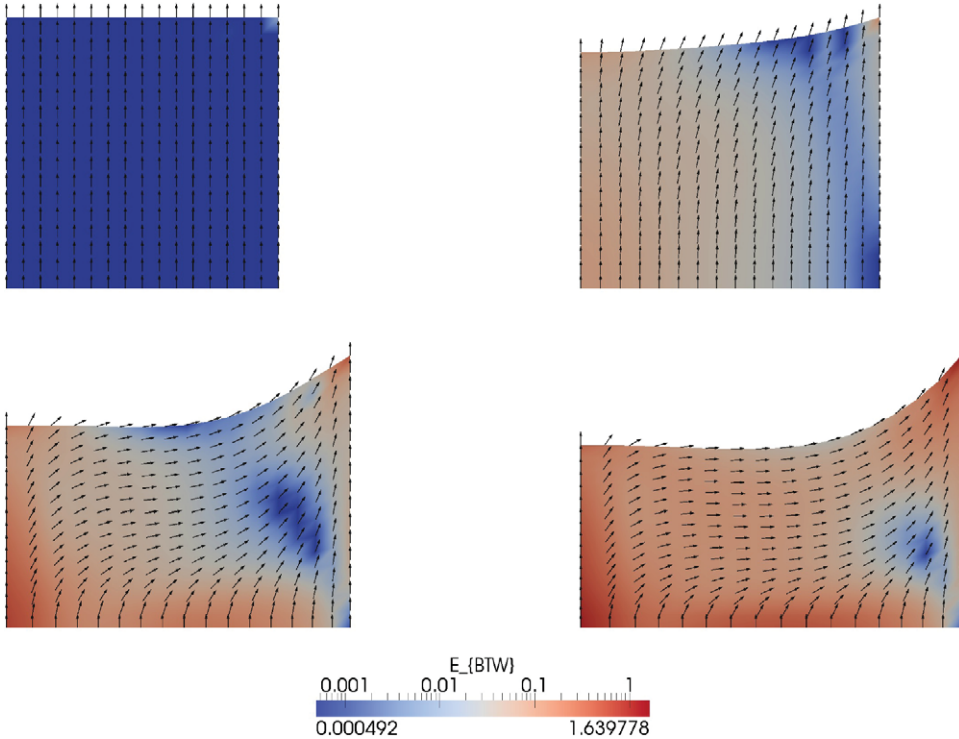


FIGURE 4. (Colour online) The solutions for AR = 1 with mesh size  $N = 16$ . From top-left to bottom-right, the strains are 0.040, 0.100, 0.264, 0.400. The domain is coloured by the BTW energy where blue corresponds to low BTW energy, while red corresponds to high BTW energy.

those strain values when  $|\kappa|$  is first and last bigger than 1, where  $\kappa$  is the curvature of the stress–strain curve. The curvature  $\kappa$  is calculated by

$$\kappa = \frac{f''}{(1 + f'^2)^{3/2}},$$

where  $f'$  is the first derivative approximated using forward difference, while  $f''$  is the second derivative approximated using central difference. From Table 3, we can see that as the mesh refines, the soft regime for AR = 1 converges to [0.076, 0.264], while the soft regime for AR = 3 converges to [0.036, 0.292].

To see what the solutions in different regimes of the stress–strain curve look like, we plot some typical solutions in Figures 4 and 5. In both figures, the top-left is a solution

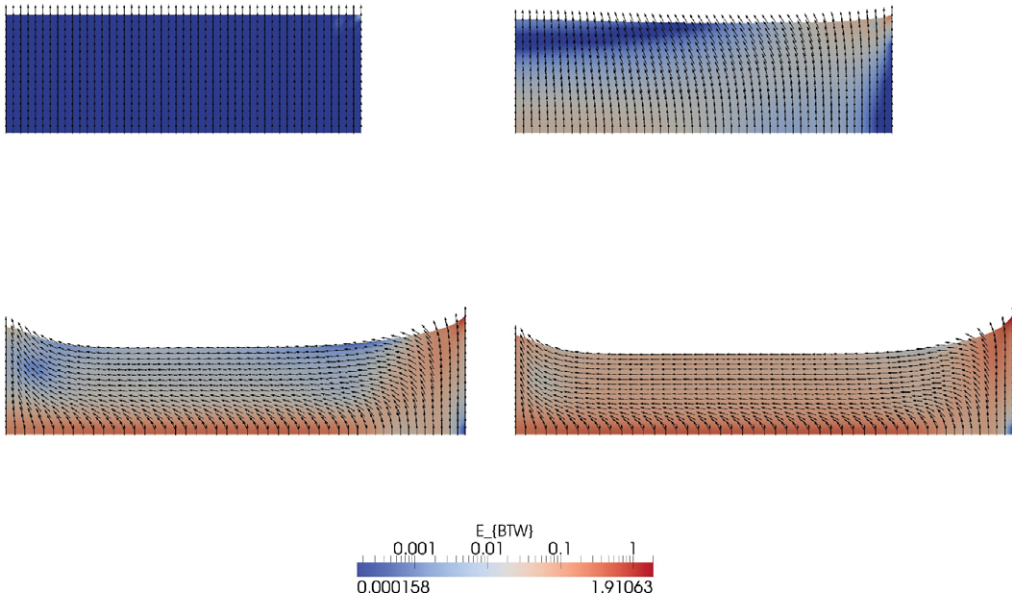


FIGURE 5. (Colour online) The solutions for  $AR = 3$  with mesh size  $N = 16$ . From top-left to bottom-right, the strains are 0.020, 0.060, 0.292, 0.400. The domain is coloured by the BTW energy where blue corresponds to low BTW energy, while red corresponds to high BTW energy.

in the first hard regime, the top-right is a solution at the start of the soft regime, the bottom-left is a solution at the end of the soft regime, while the bottom-right is a solution in the second hard regime. We can see that the solutions in the first hard regime have most directors vertical, the solutions in the second hard regime have most directors horizontal, while the solutions in the soft regime have directors rotating from vertical to horizontal. This suggests that the soft regime in the stress–strain curve might be related to the rotating of the directors. Also, we can see that the solutions in the soft regime maintain relative low BTW energy, while the solutions in the second hard regime have much higher BTW energy.

Finally we see from Figures 4 and 5 that stripe domain is not observed in these solutions. Instead, the solutions look very smooth. This might be due to the relatively coarse meshes that we have used. Due to the intrinsic high DOF of our model, the finest mesh we have used has  $N = 64$ , which might still be too coarse for the development of stripe domains. Another possible reason is that  $b = 0.0015$  might be relatively too large, which penalises the changing in  $\mathbf{n}$ , thus prevents the formation of stripe domains.

## 6 Conclusion

In this paper, we modelled the LCE using 2D BTW energy and one-constant Oseen–Frank energy. We imposed the constraint of incompressibility of the bulk, and the unity of the directors to the admissible set of  $(\mathbf{u}, \mathbf{n})$ . We proved the existence of minimiser for this energy minimisation problem. Then we converted the constrained minimisation problem to an unconstrained minimisation problem, by introducing Lagrange multipliers  $p$  and  $\lambda$ .

Next we derived the equilibrium equation and its linearisation. We reduced the linearised system to a standard saddle point system, and verified the well-posedness for some simple cases.

Next we proposed the corresponding discrete problem, which used the  $\mathbf{P}_2 \times P_1$  Taylor–Hood element for the  $(\mathbf{u}, p)$  combination and the  $\mathbf{P}_1 \times P_1$  element for the  $(\mathbf{n}, \lambda)$  combination. We also imposed the constraints that the  $L^2$  projection of  $\det(F) - 1$  is zero and that  $|\mathbf{n}| - 1$  is zero at all the mesh nodes. We proved the existence of minimiser for this discrete energy minimisation problem. Similar to the continuous case, we then introduced the Lagrange multipliers  $p$  and  $\lambda$ , derived the equilibrium equation and its linearisation, and reduced the linearised system to a standard saddle point system. We verified the well-posedness for some simple cases. For the general cases, we explained how to reduce the verification of inf-sup conditions to the computation of smallest singular value of certain matrices. Next we proved the existence and uniqueness of the Lagrange multipliers  $p$  and  $\lambda$ , under the condition that both inf-sup conditions are satisfied.

Finally we used finite element method on our model to simulate the clamped-pulling experiment, for elastomer samples with aspect ratio  $AR = 1$  or  $3$ . The orders of convergence and inf-sup values were listed. The stress–strain curves were plotted. For both  $AR = 1$  and  $3$ , the semi-soft elasticity was observed. However the stripe domain phenomenon was not observed, which might due to the relative coarse meshes in the computation and the relative large Oseen–Frank coefficient  $b = 0.0015$ .

## 7 Discussion

Although we have successfully recovered the semi-soft elasticity phenomenon, the exclusion of the stripe domain phenomenon is tentative. This is mainly because we only applied computation for meshes with size up to  $N = 64$ . That is the ratio of the edge length of the triangular elements to the edge length of the domain is around  $10^{-2}$ . However in the experiment of Finkelmann *et al.* [25, 35], the ratio of the width of the stripe domains to the edge length of the domain was around  $10^{-3}$ . Thus our mesh might be too coarse to resolve the stripe domains. We did not use finer mesh because the computational cost was already very high. Even for the mesh  $N = 64$ , we have 50,182 DOF to solve in the case  $AR = 1$ , and 149,510 DOF to solve in the case  $AR = 3$ . Another possible reason is that the Oseen–Frank coefficient  $b = 0.0015$  was too large. The zig-zag pattern of the stripe domain phenomenon naturally has very rapid change of  $\mathbf{n}$  across the domain. However a relatively large  $b$  value penalises such rapid change, and suppresses the occurrence of stripe domains. We did not take much smaller  $b$  values than 0.0015, because that would require much finer mesh and much smaller  $\Delta t$  to stabilise, which was computationally too demanding.

Stripe domain phenomenon might still be observable with our current model, if we try some more sophisticated numerical techniques. As remarked above, the main obstacle might be the computational cost. One way to get around this obstacle is to use adaptive mesh refinement. Since the stripe domains only occur in part of the elastomer domain, while  $\mathbf{n}$  in the rest of the domain is quite smooth, we can save computational cost by refining the mesh only on part of the domain. Another way to reduce the computational cost is to replace  $\mathbf{n}$  by  $(\cos(\theta), \sin(\theta))^T$ , where  $\theta$  is the azimuthal angle of the director. This



is perfectly fine because our model is in 2D. In this way, we can reduce a 2D variable to a 1D variable, and also eliminate the need to use the Lagrange multiplier  $\lambda$ .

Another direction is to replace the Oseen–Frank model by more advanced models such as Ericksen model [19] or Landau–de Gennes model [15]. Oseen–Frank energy only allows point defects, while Ericksen and Landau–de Gennes model allow line and surface defects, as well [30]. The stripe domains might have line or surface defects in the transition area between the stripes, thus using Ericksen or Landau–de Gennes model might have a better chance of capturing the stripe domain phenomenon.

### Acknowledgements

This publication was based on work supported in part by Award No. KUK-C1-013-04, made by King Abdullah University of Science and Technology (KAUST). This publication was partially supported by the National Science Foundation, Grant numbers: DMS-FRG-0456232 and DMS 1009181.

### References

- [1] ADAMS, R. & FOURNIER, J. (1975) *Sobolev Spaces*, Vol. 65, Academic Press, New York.
- [2] BALL, J. (1976) Convexity conditions and existence theorems in nonlinear elasticity. *Arch. Ration. Mech. Anal.* **63**(4), 337–403.
- [3] BLADON, P., TERENTJEV, E. & WARNER, M. (1993) Transitions and instabilities in liquid crystal elastomers. *Phys. Rev. E* **47**(6), 3838–3840.
- [4] BOCHKANOV, S. & BYSTRITSKY, V. ALGLIB [online]. URL: <http://www.alglib.net/>
- [5] BREZZI, F. & FORTIN, M. (1991) *Mixed and Hybrid Finite Element Methods*, Springer, Berlin.
- [6] CALDERER, M., LIU, C. & YAN, B. (2006) A mathematical theory for nematic elastomers with non-uniform prolate spheroids. *Advances in applied and computational mathematics*, 245–259.
- [7] CESANA, P. & DESIMONE, A. (2009) Strain-order coupling in nematic elastomers: equilibrium configurations. *Math. Models Methods Appl. Sci.* **19**, 601–630.
- [8] CHAPPELLE, D. & BATHE, K. (1993) The inf-sup test. *Comput. Struct.* **47**, 537–537.
- [9] CIARLET, P. (1988) *Mathematical Elasticity*, Vol 1, North-Holland, The Netherlands.
- [10] CLARKE, F. (1975) Generalized gradients and applications. *Trans. Am. Math. Soc.* **205**, 247–262.
- [11] CLARKE, F. (1976) A new approach to Lagrange multipliers. *Math. Oper. Res.* **1**(2), 165–174.
- [12] CLARKE, S., HOTTA, A., TAJBAKSHI, A. & TERENTJEV, E. (2001) Effect of crosslinker geometry on equilibrium thermal and mechanical properties of nematic elastomers. *Phys. Rev. E* **64**(6), 61702.
- [13] CONTI, S., DESIMONE, A. & DOLZMANN, G. (2002) Semisoft elasticity and director reorientation in stretched sheets of nematic elastomers. *Phys. Rev. E* **66**(6), 061710.
- [14] CONTI, S., DESIMONE, A. & DOLZMANN, G. (2002) Soft elastic response of stretched sheets of nematic elastomers: a numerical study. *J. the Mech. Phys. Solids* **50**(7), 1431–1451.
- [15] DE GENNES, P. & PROST, J. (1995) *The Physics of Liquid Crystals*, Oxford University Press, USA.
- [16] DESIMONE, A. (1999) Energetics of fine domain structures. *Ferroelectrics* **222**(1–4), 533–542.
- [17] DESIMONE, A. & DOLZMANN, G. (2000) Material instabilities in nematic elastomers. *Phys. D: Nonlinear Phenom.* **136**(1–2), 175–191.
- [18] DESIMONE, A. & TERESI, L. (2009) Elastic energies for nematic elastomers. *Eur. Phys. J. E: Soft Matter and Biol. Phys.* **29**(2), 191–204.

- [19] ERICKSEN, J. (1991) Liquid crystals with variable degree of orientation. *Arch. Ration. Mech. Anal.* **113**(2), 97–120.
- [20] EVANS, L. (1998) *Partial Differential Equations*. American Mathematical Society, USA.
- [21] FRANK, F. (1958) I. Liquid crystals. On the theory of liquid crystals. *Discuss. Faraday Soc.* **25**, 19–28.
- [22] HARDT, R., KINDERLEHRER, D. & LIN, F. (1986) Existence and partial regularity of static liquid crystal configurations. *Commun. Math. Phys.* **105**(4), 547–570.
- [23] HÉBERT, M., KANT, R. & DE GENNES, P. (1997) Dynamics and thermodynamics of artificial muscles based on nematic gels. *J. Physique* **7**(7), 909–919.
- [24] Hu, Q., Tai, X. & Winther, R. (2009) A saddle point approach to the computation of harmonic maps. *SIAM J. Numer. Anal.* **47**(2), 1500–1523.
- [25] KUNDLER, I. & FINKELMANN, H. (1995) Strain-induced director reorientation in nematic liquid single crystal elastomers. *Macromol. Rapid Commun.* **16**(9), 679–686.
- [26] KÜPPER, J. & FINKELMANN, H. (1994) Liquid crystal elastomers: Influence of the orientational distribution of the crosslinks on the phase behaviour and reorientation processes. *Macromol. Chem. Phys.* **195**(4), 1353–1367.
- [27] LE TALLEC, P. (1981) Compatibility condition and existence results in discrete finite incompressible elasticity. *Comput. Methods Appl. Mech. Eng.* **27**(2), 239–259.
- [28] LOGG A & WELLS, G. (2010) DOLFIN: Automated finite element computing. *ACM Trans. Math. Softw.* **37**(2), 1–28. URL: <https://launchpad.net/fenics>
- [29] LUO, C. (2010) *Modeling, Analysis and Numerical Simulation of Liquid Crystal Elastomer*, PhD Thesis, University of Minnesota, USA.
- [30] MAJUMDAR, A. & ZARNESCU, A. (2010) Landau–De Gennes theory of nematic liquid crystals: The Oseen–Frank limit and beyond. *Arch. Ration. Mech. Anal.* **196**(1), 227–280.
- [31] RUDIN, W. (1991) Functional analysis. In: *International Series in Pure and Applied Mathematics*, McGraw-Hill, USA.
- [32] TALLEC, P. L. (1994) Numerical methods for nonlinear three-dimensional elasticity. In: *Handbook of Numerical Analysis*, Vol. 3, pp. 465–622. North-Holland, The Netherlands.
- [33] VERWEY, G., WARNER, M. & TERENTJEV, E. (1996) Elastic instability and stripe domains in liquid crystalline elastomers. *J. Phys. II France* **6**, 1273–1290.
- [34] WARNER, M. & TERENTJEV, E. (2007) *Liquid Crystal Elastomers*, Oxford University Press, USA.
- [35] ZUBAREV, E., KUPTSOV, S., YURANOVA, T., TALROZE, R. & FINKELMANN, H. (1999) Monodomain liquid crystalline networks: Reorientation mechanism from uniform to stripe domains. *Liq. Cryst.* **26**(10), 1531–1540.