

# On the Analysis of the Discretized Kohn-Sham Density Functional Theory

Xin Liu\*    Zaiwen Wen†    Xiao Wang‡    Michael Ulbrich§    Yaxiang Yuan¶

April 27, 2018

**Abstract.** In this paper, we study a few theoretical issues in the discretized Kohn-Sham (KS) density functional theory (DFT). The equivalence between either a local or global minimizer of the KS total energy minimization problem and the solution to the KS equation is established under certain assumptions. The nonzero charge densities of a strong local minimizer are shown to be bounded below by a positive constant uniformly. We analyze the self-consistent field (SCF) iteration by formulating the KS equation as a fixed point map with respect to the potential. The Jacobian of these fixed point maps is derived explicitly. Both global and local convergence of the simple mixing scheme can be established if the gap between the occupied states and unoccupied states is sufficiently large. This assumption can be relaxed if the charge density is computed using the Fermi-Dirac distribution and it is not required if there is no exchange correlation functional in the total energy functional. Although our assumption on the gap is very stringent and is almost never satisfied in reality, our analysis is still valuable for a better understanding of the KS minimization problem, the KS equation and the SCF iteration.

**Key words.** Kohn-Sham total energy minimization, Kohn-Sham equation, self-consistent field iteration, nonlinear eigenvalue problem

**AMS subject classifications.** 15A18, 65F15, 47J10, 90C30

## 1 Introduction

The Kohn-Sham density functional theory in electronic structure calculations can be formulated as either a total energy minimization problem or a nonlinear eigenvalue problem. Using a suitable discretization scheme whose spatial degree of freedom is  $n$ , the electron wave functions of  $p$  occupied states can be approximated by a matrix  $X = [x_1, \dots, x_p] \in \mathbb{R}^{n \times p}$ . The charge density of electrons associated with the occupied states is defined as

$$\rho(X) := \text{diag}(XX^T), \quad (1)$$

where  $\text{diag}(A)$  denotes the vector containing the diagonal elements of the matrix  $A$ . Let  $\text{tr}(A)$  be the trace of  $A \in \mathbb{R}^{n \times n}$ , i.e., the sum of the diagonal elements of  $A$ . A commonly used discretized KS total energy function has the form of

$$E(X) := \frac{1}{4} \text{tr}(X^T LX) + \frac{1}{2} \text{tr}(X^T V_{ion} X) + \frac{1}{4} \rho^\top L^\dagger \rho + \frac{1}{2} e^\top \epsilon_{xc}(\rho), \quad (2)$$

\*State Key Laboratory of Scientific and Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, CHINA (liuxin@lsec.cc.ac.cn). Research supported in part by NSFC grants 11101409, 11331012 and 91330115, and the National Center for Mathematics and Interdisciplinary Sciences, CAS.

†Beijing International Center for Mathematical Research, Peking University, CHINA (wenzw@math.pku.edu.cn). Research supported in part by NSFC grants 11101274, 11322109 and 91330202.

‡School of Mathematical Sciences, University of Chinese Academy of Sciences, CHINA (wangxiao@ucas.ac.cn). Research supported in part by Postdoc Grant 119103S175, UCAS president grant Y35101AY00, and NSFC grant 11301505.

§Chair of Mathematical Optimization, Department of Mathematics, Technische Universität München, Boltzmannstr. 3, 85747 Garching b. München, Germany. (mulbrich@ma.tum.de).

¶State Key Laboratory of Scientific and Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, CHINA (yyx@lsec.cc.ac.cn). Research supported in part by NSFC grant 11331012.

where  $L$  is a finite dimensional representation of the Laplacian operator,  $V_{ion}$  is the ionic pseudopotentials sampled on a suitably chosen Cartesian grid,  $L^\dagger$  corresponds to the pseudo-inverse of  $L$ ,  $e$  is the column vector of all ones and  $\epsilon_{xc}(\rho)$  denotes the exchange correlation energy functional. The four terms in  $E(X)$  describe the kinetic energy, local ionic potential energy, Hartree potential energy and exchange correlation energy, respectively.

The KS total energy minimization problem solves

$$\begin{aligned} \min_{X \in \mathbb{R}^{n \times p}} \quad & E(X) \\ \text{s.t.} \quad & X^T X = I. \end{aligned} \quad (3)$$

The orthogonality constraints are imposed since the wave functions  $X$  must be orthogonal to each other due to physical constraints. It can be verified that the gradient of  $E(X)$  with respect to  $X$  is  $\nabla E(X) = H(X)X$ , where the Hamiltonian  $H(X) \in \mathbb{R}^{n \times n}$  is a matrix function

$$H(X) := \frac{1}{2}L + V_{ion} + \text{Diag}(L^\dagger \rho) + \text{Diag}(\mu_{xc}(\rho)^T e), \quad (4)$$

where  $\mu_{xc}(\rho) = \frac{\partial \epsilon_{xc}}{\partial \rho} \in \mathbb{R}^{n \times n}$  and  $\text{Diag}(x)$  denotes a diagonal matrix with  $x$  on its diagonal. The so-called KS equation is

$$\begin{aligned} H(X)X &= X\Lambda, \\ X^T X &= I, \end{aligned} \quad (5)$$

where  $\Lambda$  is a diagonal matrix consisting of  $p$  smallest eigenvalues of  $H(X)$ . The KS equation (5) is closely related to the first-order optimality conditions for (3) which are the same as (5) except that the diagonal matrix  $\Lambda$  consists of any  $p$  eigenvalues of  $H(X)$  rather than the  $p$  smallest ones.

In this paper, we first study the relationship between the KS total energy minimization problem (3) and the KS equation (5) under certain conditions. A simple counter example is provided to demonstrate that the solutions of these two problems are not necessarily the same. The second-order optimality conditions of (3) are examined based on the assumption of the existence of the second-order derivative of the exchange correlation functional [16, 29]. For a specialized exchange correlation functional, we prove that a global solution of (3) is a solution of (5) if the gap between the  $p$ th and  $(p + 1)$ st eigenvalues of the Hamiltonian  $H(X)$  is sufficiently large. The equivalence between a local minimizer of (3) and the solution (5) needs an additional assumption that the corresponding charge densities are all positive. For a strong local minimizer  $X^*$  which is defined based on the second-order sufficient optimality conditions of (3), we show that the nonzero charge densities at  $X^*$  are bounded below by a positive constant uniformly.

Our second purpose is the analysis of the most widely used approach, the self-consistent field (SCF) iteration, for solving the KS equation (5). The SCF iteration is based on computing a sequence of linear eigenvalue problems iteratively. It is well known that the basic version of SCF iteration often converges slowly or fails to converge [18] even with the help of various heuristics. A convergence analysis of the SCF iteration for solving the Hartree-Fock equations according to the optimal damping algorithm (ODA) is established in [6] and an analysis of gradient-based algorithms for the Hartree-Fock equations is proposed in [21] using Lojasiewicz inequality. The interested reader is referred to [2, 3, 4, 5, 7, 8, 9, 10, 12, 13, 20, 26] for discussion on ODA, the gradient-based algorithms and numerical analysis of DFT. A condition is identified in [30] such that the SCF iteration is a contractive fixed point iteration under a specific form of the Hamiltonian without involving any exchange correlation term. Global and local convergence of the SCF iteration for general Kohn-Sham DFT is established in [24] from an optimization point of view. Their assumptions include that the second-order derivative of the exchange correlation energy functional is uniformly bounded from above and the gap between the  $p$ th and  $(p + 1)$ st eigenvalues of the Hamiltonian  $H(X)$  is sufficiently large.

We improve the convergence results of the SCF iteration from the following three perspectives. (i) The KS equation (5) is formulated as a nonlinear system of equations (fixed point maps) respect to either the charge density or potential. Applying the differentiability of spectral operators, the Jacobian of these fixed point map is derived explicitly and analyzed. (ii) Global convergence (i.e., convergence to a stationary point from any initial solution) of the simple mixing scheme can be established when there exists a gap between  $p$ th and  $(p + 1)$ st eigenvalues of the Hamiltonian  $H(X)$ . This assumption can be relaxed for local convergence analysis, i.e., convergence behavior if the initial point is selected in a neighborhood sufficiently close to the solution of (5). If the charge density is computed using the Fermi-Dirac distribution, the assumption on the gap is not needed as long as a suitable step size for simple mixing is chosen. Our results requires much weaker conditions than the previous analysis in [24]. (iii) We propose two approximate Newton methods according to the structure of the Jacobian of the fixed point maps. The second type of our approaches is exactly the method of elliptic preconditioner proposed in [23]. Preliminary convergence results are also established for them. Although our assumption on the gap between eigenvalues of the Hamiltonian in the above three perspectives is very stringent and is almost never satisfied in reality, our analysis is still valuable for a better understanding of the KS equation and the SCF iteration.

The rest of this paper is organized as follows. A counter example between the equivalence of the KS minimization and KS equation is presented in subsection 2.1. The optimality conditions of the KS minimization problem under smoothness assumptions on the exchange functional is provided in subsection 2.2. The necessary conditions for the equivalence between a local minimizer of the KS minimization and the KS equation is established in subsection 2.3. The corresponding analysis for a global minimizer is established in subsection 2.4. Lower bounds for the charge density at local minimizers are presented in subsection 2.5. In subsection 3.1, we view the KS equation as fixed point maps with respect to the charge density or potential. The Jacobian of these fixed point maps is presented in subsection 3.2. In section 4, we establish both local and global convergence for the SCF iteration with simple mixing schemes. Two approximate Newton approaches and their convergence properties are discussed in section 5.

## 2 Equivalence Between the KS Total Energy Minimization and the KS Equation

### 2.1 A Counter Example

The following three-dimensional toy example shows that a solution of the KS equation is not necessary a global optimal solution of the KS total energy minimization problem. Let  $n = 3$ ,  $p = 1$  and choose

$$L = \begin{pmatrix} 1.4299 & -0.2839 & -0.4056 \\ -0.2839 & 1.1874 & 0.2678 \\ -0.4056 & 0.2678 & 1.3826 \end{pmatrix}, \quad V_{ion} = 0, \text{ and } \epsilon_{xc}(\rho) = 0.$$

It can be verified numerically that  $X^* = (0.3683 \quad -0.6188 \quad 0.6939)^T$  is a global minimizer of (3). On the other hand, we have

$$H(X^*) = \frac{1}{2}L + \text{Diag}(L^\dagger \rho(X^*)) = \begin{pmatrix} 0.9735 & -0.1419 & -0.2028 \\ -0.1419 & 0.8955 & 0.1339 \\ -0.2028 & 0.1339 & 1.0569 \end{pmatrix},$$

and  $X^*$  is an eigenvector associated with the second smallest eigenvalue of  $H(X^*)$ . Therefore, the equivalence between the KS total energy minimization and the KS equation only holds under certain assumptions. For this counter

example, our assumptions in subsections 2.3 and 2.4 do not hold because the gap between the eigenvalues of  $H(X^*)$  is  $\delta = 0.046$  and it is smaller than  $\|L^\dagger\|_2 = 1$ . We should point out that the above example may not exist in the practice of DFT.

## 2.2 Optimality Conditions Under Smoothness Assumptions on $\epsilon_{xc}(\rho)$

The Lagrangian function of the minimization problem (3) is

$$\mathcal{L}(X, \Lambda) := E(X) - \frac{1}{2} \text{tr}(\Lambda(X^T X - I)).$$

Suppose  $X$  is a local minimizer of (3). It follows from  $X^T X = I$  that the linear independence constraint qualification is satisfied. Hence, there exists a Lagrange multiplier  $\Lambda$  such that the first-order optimality conditions hold:

$$\nabla_X \mathcal{L}(X, \Lambda) = H(X)X - X\Lambda = 0 \text{ and } X^T X = I. \quad (6)$$

Multiplying both sides of the first equality in (6) by  $X^T$  and using  $X^T X = I$ , we have  $\Lambda = X^T H(X)X$ , which is a symmetric matrix. Note that  $E(XQ) = E(X)$  and  $H(XQ) = H(X)$  hold for any orthogonal matrix  $Q \in \mathbb{R}^{p \times p}$ . Hence, if  $X$  is a stationary point, any matrix in the set  $\{XQ \mid Q \in \mathbb{R}^{p \times p} \text{ and } Q^T Q = I\}$  is also a stationary point, and their objective values are the same. Let  $\tilde{Q}\tilde{\Lambda}\tilde{Q}^T$  be the eigenvalue decomposition of  $X^T H(X)X$  and  $\tilde{X} := X\tilde{Q}$ . Then the Lagrangian multiplier  $\tilde{\Lambda} = \tilde{X}^T H(\tilde{X})\tilde{X}$  is a diagonal matrix whose entries are the eigenvalues of  $H(X)$ .

Let  $\mathcal{L}(\mathbb{R}^{n \times p}, \mathbb{R}^{n \times p})$  denote the space of linear operators which map  $\mathbb{R}^{n \times p}$  to  $\mathbb{R}^{n \times p}$ . The Fréchet derivative of  $\nabla E(X)$  is defined as the (unique) function  $\nabla^2 E : \mathbb{R}^{n \times p} \rightarrow \mathcal{L}(\mathbb{R}^{n \times p}, \mathbb{R}^{n \times p})$  such that

$$\lim_{\|S\|_F \rightarrow 0} \frac{\|\nabla E(X+S) - \nabla E(X) - \nabla^2 E(X)(S)\|_F}{\|S\|_F} = 0.$$

The next lemma shows an explicit form of the Hessian operator, if the exchange correlation energy is second-order differentiable.

**Lemma 2.1** (Lemma 2.1 in [29]). *Suppose that  $\epsilon_{xc}(\rho(X))$  is twice differentiable with respect to  $\rho(X)$ . Given a direction  $S \in \mathbb{R}^{n \times p}$ , the Hessian-vector product of  $E(X)$  is*

$$\nabla^2 E(X)[S] = H(X)S + 2\text{Diag}(J(\rho)\text{diag}(SX^T))X, \quad (7)$$

where

$$J(\rho) := L^\dagger + \partial\mu_{xc}(\rho)e. \quad (8)$$

Consequently, the second-order necessary and sufficient optimality conditions can be obtained from Theorems 12.5 and 12.6 in [25], respectively.

**Theorem 2.2.** *1) Suppose that  $X$  is a local minimizer of problem (3) and  $\epsilon_{xc}(\rho(X))$  is twice differentiable with respect to  $\rho(X)$ . Then, for all  $S \in \mathcal{T}(X)$ , it holds*

$$\text{tr}(S^T H(X)S - \Lambda S^T S) + 2\text{diag}(XS^T)^T J \text{diag}(XS^T) \geq 0, \quad (9)$$

where  $\Lambda = X^T H(X)X$  and

$$\mathcal{T}(X) := \{S \mid X^T S + S^T X = 0\}. \quad (10)$$

2) Suppose that  $X \in \mathbb{R}^{n \times p}$  satisfies (6) with a symmetric matrix  $\Lambda$  and (9) holds with a strict inequality for all  $0 \neq S \in \mathcal{T}(X)$ . Then  $X$  is a strict local minimizer for problem (3).

*Proof.* It follows from Theorem 12.5 in [25] that the second-order necessary condition for  $X$  to be a local minimizer of (3) is

$$\langle S, \nabla_{XX}^2 \mathcal{L}(X, \Lambda)[S] \rangle \geq 0, \quad \text{for all } S \in \mathcal{T}(X). \quad (11)$$

Using Lemma 2.1 and the fact that

$$\text{tr}(X^T \text{Diag}(y)Z) = y^T \text{diag}(ZX^T), \quad \text{for all } X, Z \in \mathbb{R}^{n \times p}, y \in \mathbb{R}^n,$$

we obtain

$$\begin{aligned} \langle S, \nabla_{XX}^2 \mathcal{L}(X, \Lambda)[S] \rangle &= \text{tr}(S^T \nabla^2 E(X)[S] - \Lambda S^T S) \\ &= \text{tr}(S^T H(X)S + 2S^T \text{Diag}(J \text{diag}(SX^T))X - \Lambda S^T S) \\ &= \text{tr}(S^T H(X)S - \Lambda S^T S) + 2\text{diag}(XS^T)^T J \text{diag}(XS^T), \end{aligned}$$

which together with (11) yields (9). The second part is a direct application of Theorem 12.6 in [25].  $\square$

An equivalent formulation of the tangent space (10) is

$$\mathcal{T}(X) = \{S := XK + \mathbf{P}_X^\perp Z \mid K = -K^T \in \mathbb{R}^{p \times p}, Z \in \mathbb{R}^{n \times p}\}, \quad (12)$$

where  $\mathbf{P}_X^\perp := I - XX^T$ . Hence, the second-order optimality conditions in Theorem 2.2 can be presented in terms of an arbitrary  $Z \in \mathbb{R}^{n \times p}$  similar to the analysis of maximization of the sum of the trace ratio on the Stiefel Manifold in [31].

**Theorem 2.3.** 1) Suppose that  $X$  is a local minimizer of problem (3) and  $\epsilon_{xc}(\rho(X))$  is twice differentiable with respect to  $\rho(X)$ . Then for all  $Z \in \mathbb{R}^{n \times p}$ , it holds

$$\begin{aligned} &\text{tr}(Z^T H(X)Z) + \text{tr}(X^T Z \Lambda Z^T X) - \text{tr}(Z^T X \Lambda X^T Z) - \text{tr}(Z \Lambda Z^T) \\ &+ 2\text{diag}(XZ^T \mathbf{P}_X^\perp)^T J \text{diag}(XZ^T \mathbf{P}_X^\perp) \geq 0. \end{aligned} \quad (13)$$

2) Suppose that  $X \in \mathbb{R}^{n \times p}$  satisfies (6) with a symmetric matrix  $\Lambda$  and (13) holds with a strict inequality for all  $\mathbf{P}_X^\perp Z \neq 0$ . Then  $X$  is a strict local minimizer for problem (3).

*Proof.* Using (6) and the definition of  $\mathbf{P}_X^\perp$ , we obtain  $\mathbf{P}_X^\perp \mathbf{P}_X^\perp = \mathbf{P}_X^\perp$ ,  $\mathbf{P}_X^\perp X = 0$  and  $\mathbf{P}_X^\perp H(X)X = 0$ . For any  $S = XK + \mathbf{P}_X^\perp Z$ , it holds

$$\begin{aligned} \text{tr}(S^T H(X)S) &= \text{tr}(K^T X^T H(X)XK) + \text{tr}(Z^T \mathbf{P}_X^\perp H(X) \mathbf{P}_X^\perp Z) \\ &= \text{tr}(K^T \Lambda K) + \text{tr}(Z^T H(X)Z) - \text{tr}(Z^T H(X)X X^T Z) \\ &= \text{tr}(K^T \Lambda K) + \text{tr}(Z^T H(X)Z) - \text{tr}(Z^T X \Lambda X^T Z). \end{aligned} \quad (14)$$

It can be verified that  $S^T S = K^T K + Z^T \mathbf{P}_X^\perp Z$ , which yields

$$\text{tr}(\Lambda S^T S) = \text{tr}(K^T K \Lambda) + \text{tr}(Z^T Z \Lambda) - \text{tr}(Z^T X X^T Z \Lambda)$$

$$= \operatorname{tr}(K^T \Lambda K) + \operatorname{tr}(Z \Lambda Z^T) - \operatorname{tr}(X^T Z \Lambda Z^T X), \quad (15)$$

where the last equality holds because of  $K = -K^T$ . Since it holds

$$\operatorname{diag}(X K^T X^T) = \frac{1}{2}(\operatorname{diag}(X K^T X^T) + \operatorname{diag}(X K X^T)) = \frac{1}{2}\operatorname{diag}(X(K + K^T)X^T) = 0,$$

we obtain

$$\operatorname{diag}(X S^T) = \operatorname{diag}(X K^T X^T) + \operatorname{diag}(X Z^T \mathbf{P}_X^\perp) = \operatorname{diag}(X Z^T \mathbf{P}_X^\perp),$$

which together with (14) and (15) gives (13). The proof of the second part follows directly from Theorem 2.2.  $\square$

### 2.3 Necessary Condition for Local Minimizers

In this subsection, we establish a necessary condition under which a local minimizer of (3) is a solution of a modification of the KS equation (5). Our discussion is restricted to a special exchange correlation functional

$$e^T \epsilon_{xc}(\rho) = -\frac{3}{4}\gamma \rho^T \rho^{\frac{1}{3}}, \quad (16)$$

where  $\gamma = 2\left(\frac{3}{\pi}\right)^{1/3}$  and  $\rho^{\frac{1}{3}}$  denotes the component-wise cubic root of the vector  $\rho$ . The next result shows that the charge density  $\rho$  is bounded.

**Lemma 2.4.** *Let  $X \in \mathbb{R}^{n \times p}$  satisfy  $X^T X = I$ , and  $\rho$  be defined by (1). We have*

$$0 \leq \rho_i \leq 1, \text{ for all } i = 1, \dots, n. \quad (17)$$

*Proof.* The inequality (17) holds from  $X^T X = I$  and the fact that  $\rho_i = \sum_{j=1}^p X_{ij}^2$  for all  $i = 1, \dots, n$ .  $\square$

Our analysis relies on the gap between the  $p$ th and  $(p+1)$ st eigenvalues of  $H(X)$ .

**Assumption 2.5.** *Let  $\lambda_1 \leq \dots \leq \lambda_p \leq \lambda_{p+1} \leq \dots \leq \lambda_n$  be the eigenvalues of a given symmetric matrix  $H \in \mathbb{R}^{n \times n}$ . There exists a positive constant  $\delta$  such that  $\lambda_{p+1} - \lambda_p \geq \delta$ .*

Note that  $E(X)$  may not be second-order differentiable since some components  $\rho_i(X)$  can be zero. Let  $\mathcal{I}$  be the collection of indices of the nonzero components of  $\rho(X)$ , i.e.,

$$\mathcal{I} = \{i \mid \rho_i(X) \neq 0, i = 1, \dots, n\}. \quad (18)$$

Then the complement set  $\bar{\mathcal{I}}$  of  $\mathcal{I}$  is the set of indices of the zero components of  $\rho(X)$ . Let  $r$  be the cardinality of  $\mathcal{I}$ . We have  $r \geq p$  by the orthogonality of  $X$ . If  $\mathcal{I} = \{\alpha_1, \dots, \alpha_r\}$ , we define the submatrices  $X_{\mathcal{I}}$  and  $L_{\mathcal{I}\mathcal{I}}$  as

$$X_{\mathcal{I}} = \begin{pmatrix} X_{\alpha_1,1}, \dots, X_{\alpha_1,p} \\ \dots \\ X_{\alpha_r,1}, \dots, X_{\alpha_r,p} \end{pmatrix}, \text{ and } L_{\mathcal{I}\mathcal{I}} = \begin{pmatrix} L_{\alpha_1,1}, \dots, L_{\alpha_1,\alpha_r} \\ \dots \\ L_{\alpha_r,1}, \dots, L_{\alpha_r,\alpha_r} \end{pmatrix}.$$

The notations  $(V_{ion})_{\mathcal{I}\mathcal{I}}$ ,  $L_{\mathcal{I}\mathcal{I}}^\dagger$ ,  $H_{\mathcal{I}\mathcal{I}}(X)$  and  $\Lambda_{\mathcal{I}\mathcal{I}}$  are defined similar to  $L_{\mathcal{I}\mathcal{I}}$ .

The following theorem shows that a local minimizer  $X^*$  of the KS total energy minimization (3) is a solution of KS equation (5) if all rows of  $X^*$  are nonzero and Assumption 2.5 holds with a sufficiently large gap  $\delta$ .

**Theorem 2.6.** Suppose that  $X^*$  is a local minimizer of (3) using (16) and  $\Lambda^* = (X^*)^\top H(X^*)X^*$  is a diagonal matrix. Let  $\mathcal{I}^*$  be the index set of  $X^*$  defined as (18). If Assumption 2.5 holds at  $H(X^*)$  with a constant  $\delta$  satisfying

$$\delta > 2 \left( \|L^\dagger\|_2 - \frac{\gamma}{3} \right), \quad (19)$$

then it holds

$$\begin{aligned} H_{\mathcal{I}^* \mathcal{I}^*}(X^*) X_{\mathcal{I}^*}^* &= X_{\mathcal{I}^*}^* \Lambda^*, \\ (X_{\mathcal{I}^*}^*)^\top X_{\mathcal{I}^*}^* &= I, \end{aligned} \quad (20)$$

and the diagonal of  $\Lambda^*$  consists of the  $p$  smallest eigenvalues of  $H_{\mathcal{I}^* \mathcal{I}^*}(X^*)$ .

*Proof.* It can be verified that  $X^*$  is a local minimizer of the restricted problem

$$\begin{aligned} \min_{X \in \mathbb{R}^{n \times p}} \quad & E(X) \\ \text{s.t.} \quad & X^\top X = I, \quad X_{\mathcal{I}^*} = 0. \end{aligned} \quad (21)$$

Hence,  $X_{\mathcal{I}^*}^*$  is a local minimizer of the reduced problem

$$\begin{aligned} \min_{\hat{X} \in \mathbb{R}^{r \times p}} \quad & \hat{E}(\hat{X}) := \frac{1}{4} \text{tr}(\hat{X}^\top L_{\mathcal{I}^* \mathcal{I}^*} \hat{X}) + \frac{1}{2} \text{tr}(\hat{X}^\top (V_{ion})_{\mathcal{I}^* \mathcal{I}^*} \hat{X}) + \frac{1}{4} \rho(\hat{X})^\top L_{\mathcal{I}^* \mathcal{I}^*}^\dagger \rho(\hat{X}) - \frac{3}{4} \gamma \rho(\hat{X})^\top \rho(\hat{X})^{\frac{1}{3}}, \\ \text{s.t.} \quad & \hat{X}^\top \hat{X} = I. \end{aligned} \quad (22)$$

The structure of the energy functional  $E(X)$  implies  $\nabla \hat{E}(X_{\mathcal{I}^*}^*) = H_{\mathcal{I}^* \mathcal{I}^*}(X^*) X_{\mathcal{I}^*}^*$  and  $(X_{\mathcal{I}^*}^*)^\top H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*) X_{\mathcal{I}^*}^* = \Lambda^*$ . These facts together with the first-order optimality of (22) at  $X_{\mathcal{I}^*}^*$  yield (20).

It is obvious that the diagonal entries of  $\Lambda^*$  are the eigenvalues of  $H_{\mathcal{I}^* \mathcal{I}^*}(X^*)$ . Suppose that they are not the  $p$  smallest eigenvalues of  $H_{\mathcal{I}^* \mathcal{I}^*}(X^*)$ . For convenience, we denote the eigenvalues of  $H_{\mathcal{I}^* \mathcal{I}^*}(X^*)$  in an ascending order as  $\hat{\lambda}_1 \leq \dots \leq \hat{\lambda}_r$ , and their corresponding eigenvectors are  $u_i, i = 1, \dots, r$ , where  $r = |\mathcal{I}^*|$ . Let  $x_i, 1 \leq i \leq p$ , be the  $i$ th column for  $X_{\mathcal{I}^*}^*$ . Without loss of generality, let  $x_1$  be associated with an eigenvalue greater than  $\hat{\lambda}_p$ , and  $u_i (i \leq p)$  be an eigenvector associated with an eigenvalue less than or equal to  $\hat{\lambda}_p$  but not be a column of  $X_{\mathcal{I}^*}^*$ . The Assumption 2.5 implies that  $u_i \notin \text{span}\{X_{\mathcal{I}^*}^*\}$ . Let  $V$  be a matrix whose columns satisfy

$$v_j = \begin{cases} u_i & \text{if } j = 1, \\ x_j & \text{if } j = 2, \dots, p. \end{cases}$$

Since the function  $\hat{E}(\hat{X})$  is twice differentiable at  $X_{\mathcal{I}^*}^*$  according to the definition of  $\mathcal{I}^*$ . Therefore, an application of Theorem 2.3 gives

$$\begin{aligned} \Delta &:= \text{tr}(V^\top H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*) V) + \text{tr}((X_{\mathcal{I}^*}^*)^\top V \Lambda^* V^\top X_{\mathcal{I}^*}^*) - \text{tr}(V^\top X_{\mathcal{I}^*}^* \Lambda^* (X_{\mathcal{I}^*}^*)^\top V) - \text{tr}(V \Lambda^* V^\top) \\ &\quad + 2 \text{diag}(X_{\mathcal{I}^*}^* V^\top \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp)^\top \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} \text{Diag} \left( \rho(X_{\mathcal{I}^*}^*)^{-\frac{2}{3}} \right) \right) \text{diag}(X_{\mathcal{I}^*}^* V^\top \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp) \\ &\geq 0. \end{aligned} \quad (23)$$

It follows from that  $V$  is an orthonormal eigenbasis of  $H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*)$  and Assumption 2.5 that

$$\text{tr}(V^\top H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*) V) - \text{tr}((X_{\mathcal{I}^*}^*)^\top H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*) X_{\mathcal{I}^*}^*) \leq \hat{\lambda}_i - \hat{\lambda}_{p+1} \leq -\delta. \quad (24)$$

Since  $u_i \notin \text{span}\{X_{\mathcal{I}^*}^*\}$ , we obtain

$$(X_{\mathcal{I}^*}^*)^T V = V^T X_{\mathcal{I}^*}^* = I - e_1 e_1^T, \quad (25)$$

$$X_{\mathcal{I}^*}^* V^T \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp = x_1 u_i^T, \quad (26)$$

which further give

$$\begin{aligned} \Delta &= \text{tr}(V^T H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*) V) - \text{tr}(\Lambda^*) + 2 \text{diag}(x_1 u_i^T)^T \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} \text{Diag} \left( \rho(X_{\mathcal{I}^*}^*)^{-\frac{2}{3}} \right) \right) \text{diag}(x_1 u_i^T) \\ &\leq -\delta + 2 \max \left\{ \lambda_{\max} \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} \text{Diag} \left( \rho(X_{\mathcal{I}^*}^*)^{-\frac{2}{3}} \right) \right), 0 \right\} \\ &\leq -\delta + 2 \max \left\{ \lambda_{\max} \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} I \right), 0 \right\} \\ &\leq -\delta + 2 \max \left\{ \left( \|L_{\mathcal{I}^* \mathcal{I}^*}^\dagger\|_2 - \frac{\gamma}{3} \right), 0 \right\} \\ &< 0, \end{aligned} \quad (27)$$

where the first inequality uses (24) and the fact that  $\|\text{diag}(x_1 u_i^T)\|_2^2 \leq 1$ , the second inequality follows from  $\rho \in [0, 1]$ , the third inequality uses the fact that  $\|L_{\mathcal{I}^* \mathcal{I}^*}^\dagger\|_2 \leq \|L^\dagger\|_2$  since the largest/smallest eigenvalue of a matrix is no less/greater than the largest/smallest eigenvalue of its any principal submatrix, and the last inequality (27) is due to (19). However, (27) is a contradiction to (23). This completes the proof.  $\square$

## 2.4 Necessary Condition for Global Minimizers

In this subsection, we consider whether a global minimizer of (3) is a solution of the KS equation (5) under the exchange correlation functional (16). We first show the following inequality.

**Lemma 2.7.** *It holds for all  $a, b \in [0, 1]$  that*

$$(a - b)^2(3a^2 + 2ab + b^2) = 3a^4 - 4a^3b + b^4 \geq \frac{2}{3}(a^3 - b^3)^2.$$

*Proof.* The inequality holds for  $a = 0$  or  $b = 0$ . Consider the case on  $a \geq b > 0$ . Introducing the variable  $t = b/a \in (0, 1]$  yields

$$a^4(3 - 4t + t^4) - \frac{2}{3}a^6(1 - t^3)^2 \geq a^6 f(t),$$

where  $f(t) = 3 - 4t + t^4 - \frac{2}{3}(1 - t^3)^2$ . Since  $f'(t) = (t^3 - 1)(4 - 4t^2) \leq 0$  for all  $t \in [0, 1]$ , we have  $f(t) \geq f(1) = 0$  for all  $t \in [0, 1]$ , and then the inequality is proved. The case on  $b \geq a > 0$  can be proved in a similar fashion.  $\square$

The next theorem establishes the equivalence based on estimating the difference of total energy function values.

**Theorem 2.8.** *Suppose that  $X^*$  is a global minimizer of (3) using (16). If Assumption 2.5 holds at  $H(X^*)$  with a constant  $\delta$  satisfying*

$$\delta > p \left( \|L^\dagger\|_2 - \frac{\gamma}{3} \right), \quad (28)$$

*then  $X^*$  must be an orthonormal eigenbasis of  $H(X^*)$  corresponding to its  $p$  smallest eigenvalues, namely, a solution of the KS equation (5).*

*Proof.* Suppose that  $X^*$  is not but  $Y$  is an orthonormal eigenbasis of  $H(X^*)$  corresponding to its  $p$  smallest eigenvalues. Since  $X^*$  must be an orthonormal eigenbasis of  $H(X^*)$  and using Assumption 2.5, we have

$$\Delta H(Y, X^*) := \text{tr}(Y^T H(X^*) Y) - \text{tr}((X^*)^T H(X^*) X^*) \leq \lambda_p(H(X^*)) - \lambda_{p+1}(H(X^*)) \leq -\delta. \quad (29)$$

Applying Lemmas 2.4 and 2.7 gives

$$\sum_{i=1}^n \left( \rho(Y)_i^{\frac{1}{3}} - \rho(X^*)_i^{\frac{1}{3}} \right)^2 \left( 3\rho(Y)_i^{\frac{2}{3}} + 2\rho(Y)_i^{\frac{1}{3}} \rho(X^*)_i^{\frac{1}{3}} + \rho(X^*)_i^{\frac{2}{3}} \right) \geq \frac{2}{3} \|\rho(Y) - \rho(X^*)\|_2^2. \quad (30)$$

It follows from Lemma 2.4 that

$$\begin{aligned} \|\rho(Y) - \rho(X^*)\|^2 &\leq (1 - \rho(Y))^T \rho(X^*) + (1 - \rho(X^*))^T \rho(Y) \\ &\leq \mathbf{1}^T \rho(X^*) + \mathbf{1}^T \rho(Y) = \text{tr}(X X^T) + \text{tr}(Y Y^T) \\ &= 2p. \end{aligned} \quad (31)$$

Using the relationship  $\text{tr}(Y^T \text{Diag}(L^\dagger \rho(X^*)) Y) = \rho(Y)^T L^\dagger \rho(X^*)$ , the inequalities (29), (30) and (31), and the assumption (28), we obtain

$$\begin{aligned} \Delta E(Y, X^*) &= E(Y) - E(X^*) \\ &= \frac{1}{2} \Delta H(Y, X^*) + \frac{1}{4} (\rho(Y)^T L^\dagger \rho(Y) - \rho(X^*)^T L^\dagger \rho(X^*)) - \frac{3\gamma}{8} (\rho(Y)^T \rho(Y)^{\frac{1}{3}} - \rho(X^*)^T \rho(X^*)^{\frac{1}{3}}) \\ &\quad - \frac{1}{2} \text{tr}(Y^T \text{Diag}(L^\dagger \rho(X^*) - \gamma \rho(X^*)^{\frac{1}{3}}) Y) + \frac{1}{2} \text{tr}(X^T \text{Diag}(L^\dagger \rho(X^*) - \gamma \rho(X^*)^{\frac{1}{3}}) X^*) \\ &= \frac{1}{2} \Delta H(Y, X^*) + \frac{1}{4} (\rho(Y)^T L^\dagger \rho(Y) - \rho(X^*)^T L^\dagger \rho(X^*)) - \frac{3\gamma}{8} (\rho(Y)^T \rho(Y)^{\frac{1}{3}} - \rho(X^*)^T \rho(X^*)^{\frac{1}{3}}) \\ &\quad - \frac{1}{2} (\rho(Y)^T L^\dagger \rho(X^*) - \rho(X^*)^T L^\dagger \rho(X^*)) + \frac{1}{2} \gamma (\rho(Y)^T \rho(X^*)^{\frac{1}{3}} - \rho(X^*)^T \rho(X^*)^{\frac{1}{3}}) \\ &= \frac{1}{2} \Delta H(Y, X^*) + \frac{1}{4} (\rho(Y) - \rho(X^*))^T L^\dagger (\rho(Y) - \rho(X^*)) \\ &\quad - \frac{\gamma}{8} \sum_{i=1}^n \left( \rho(Y)_i^{\frac{1}{3}} - \rho(X^*)_i^{\frac{1}{3}} \right)^2 \left( 3\rho(Y)_i^{\frac{2}{3}} + 2\rho(Y)_i^{\frac{1}{3}} \rho(X^*)_i^{\frac{1}{3}} + \rho(X^*)_i^{\frac{2}{3}} \right) \\ &\leq -\frac{\delta}{2} + \left( \frac{\|L^\dagger\|_2}{4} - \frac{\gamma}{12} \right) \|\rho(Y) - \rho(X^*)\|_2^2 \\ &\leq -\frac{\delta}{2} + \left( \frac{\|L^\dagger\|_2}{4} - \frac{\gamma}{12} \right) (2p) \\ &< 0, \end{aligned}$$

which is a contradiction to the fact that  $X^*$  is a global minimizer. This completes the proof.  $\square$

**Remark 2.9.** When the exchange correlation function  $\epsilon_{xc}(\rho)$  is equal to zero, our condition (28) becomes  $\delta > p\|L^\dagger\|_2$ , which is much weaker than the condition  $\delta > 12p\sqrt{n}\|L^\dagger\|_2$  in Theorem 1 of [24].

## 2.5 Lower Bounds for the Charge Density of Local Minimizers

The exchange correlation energy functional is twice differentiable if all components of  $\rho(X)$  are positive. However, the second-order derivative may not be bounded at an arbitrary point  $X$ . In this subsection, we provides a few lower

bounds for the charge density at certain types local minimizers. These properties are useful for our analysis on the KS equation.

Traditionally, a point  $x^*$  is called a strong local minimizer [1, 15] of a function  $f : \mathbb{R}^n \mapsto \mathbb{R}$ , if there exists a constant  $\kappa > 0$  and a neighborhood  $U$  of  $x^*$  such that the inequality

$$f(x) \geq f(x^*) + \kappa \|x - x^*\|_2^2 \quad (32)$$

holds for any  $x \in U$ . Here, we define a strong local minimizer based on the second-order optimality conditions.

**Definition 2.10.** A point  $X^*$  is called a strong local minimizer of (3) using (16) if and only if  $X_{\mathcal{I}^*}^*$  is local minimizer of (22) and there exists a constant  $\kappa > 0$  such that, for all  $Z \in \mathbb{R}^{n \times p}$ ,

$$\begin{aligned} & \text{tr}(Z^T H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*) Z) + \text{tr}((X_{\mathcal{I}^*}^*)^T Z \Lambda^* Z^T X_{\mathcal{I}^*}^*) - \text{tr}(Z^T X_{\mathcal{I}^*}^* \Lambda^* (X_{\mathcal{I}^*}^*)^T Z) - \text{tr}(Z \Lambda^* Z^T) \\ & + 2 \text{diag}((X_{\mathcal{I}^*}^*) Z^T \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp)^T \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} \text{Diag} \left( \rho(X_{\mathcal{I}^*}^*)^{-\frac{2}{3}} \right) \right) \text{diag}((X_{\mathcal{I}^*}^*) Z^T \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp) \geq \kappa \|Z\|_F^2, \end{aligned} \quad (33)$$

where  $\Lambda^* = (X_{\mathcal{I}^*}^*)^T H_{\mathcal{I}^* \mathcal{I}^*}(X^*) X_{\mathcal{I}^*}^*$  and  $\mathcal{I}^*$  is the index set of  $X^*$  defined as (18).

Our condition (33) is weaker than (32) applying to problem (3) when the total energy  $E(X)$  is twice differentiable. The next result shows that the charge densities at a strong local minimizer are bounded below uniformly if they are positive.

**Theorem 2.11.** Suppose that  $L$  is positive semidefinite and  $X^*$  is a strong local minimizer of (3) satisfying Definition 2.10. Let

$$\bar{c} := \min\{1, c_1, \dots, c_n\} \text{ and } c_i := \min_{j \neq i} \left( \frac{\gamma}{3(L_{ii}^\dagger - 2L_{ij}^\dagger + L_{jj}^\dagger)} \right)^{\frac{3}{2}}. \quad (34)$$

Then it holds:

$$\text{for any } i \in \{1, 2, \dots, n\}, \quad \rho_i(X^*) \in [0, \bar{c}] \Rightarrow \rho_i(X^*) = 0. \quad (35)$$

*Proof.* For convenience, we denote  $\rho_{\mathcal{I}^*}^* = \rho(X_{\mathcal{I}^*}^*)$ . If there exists a row  $j$  in  $X_{\mathcal{I}^*}^*$  such that either 1 or  $-1$  is an entry of this row, then this row has only one nonzero entry according to the orthonormality of  $X_{\mathcal{I}^*}^*$ . Hence,  $(\rho_{\mathcal{I}^*}^*)_j = 1$  and (35) holds at  $j$ .

We next consider the components in the set  $\mathcal{J} := \{j \mid j \in \mathcal{I}^* \text{ and } |(X_{\mathcal{I}^*}^*)_{js}| < 1, s = 1, \dots, p\}$ . For any given  $j \in \mathcal{J}$ , there exists a nonzero entry, denoted as  $(X_{\mathcal{I}^*}^*)_{js}$ , in the  $j$ -th row of  $X_{\mathcal{I}^*}^*$ . Since  $|(X_{\mathcal{I}^*}^*)_{js}| < 1$ , there exists at least another nonzero entry, denoted as  $(X_{\mathcal{I}^*}^*)_{is}$ , in the  $s$ -th column of  $X_{\mathcal{I}^*}^*$  due to the orthonormality of  $X_{\mathcal{I}^*}^*$ . For simplicity, let  $x_l$ ,  $l = 1, \dots, p$ , be the  $l$ -th column of  $X_{\mathcal{I}^*}^*$  and set  $r = |\mathcal{I}^*|$ ,  $x_{js} = (X_{\mathcal{I}^*}^*)_{js}$  and  $x_{is} = (X_{\mathcal{I}^*}^*)_{is}$ . Define a vector  $z \in \mathbb{R}^r$  whose  $l$ -th component ( $l = 1, \dots, p$ ) is

$$z_l = \begin{cases} \frac{x_{is}}{\sqrt{x_{is}^2 + x_{js}^2}}, & \text{if } l = j; \\ \frac{-x_{js}}{\sqrt{x_{is}^2 + x_{js}^2}}, & \text{if } l = i; \\ 0, & \text{otherwise.} \end{cases} \quad (36)$$

A short calculation gives  $\|z\|_2 = 1$ ,  $z^\top x_s = 0$  and

$$\text{diag}(zx_s^\top) = \frac{x_{is}x_{js}}{\sqrt{x_{is}^2 + x_{js}^2}} e_{(j,-i)}, \quad (37)$$

where  $e_{(j,-i)} \in \mathbb{R}^r$  has 1 on its  $j$ -th entry,  $-1$  on its  $i$ -th entry and 0 elsewhere.

For  $a \in [0, 1]$ , let  $Z_a \in \mathbb{R}^{n \times p}$  be a matrix whose  $s$ -th column is  $az + \sqrt{1-a^2}x_s$  and all other columns are zero. Without loss of generality, let  $\hat{\lambda}_1 \leq \dots \leq \hat{\lambda}_r$  be the eigenvalues of  $H_{\mathcal{I}^* \mathcal{I}^*}(X^*)$  in the ascending order, and  $x_s$  be an eigenvector of  $H_{\mathcal{I}^* \mathcal{I}^*}(X^*)$  associated with  $\hat{\lambda}_s$ ,  $s \in \{1, \dots, r\}$ . Then, we obtain

$$\text{tr}(Z_a^\top H_{\mathcal{I}^* \mathcal{I}^*}(X^*) Z_a) \leq a^2 \hat{\lambda}_r + (1-a^2) \hat{\lambda}_s, \quad (38)$$

$$\text{tr}(Z_a \Lambda^* Z_a^\top) = \text{tr}(\Lambda^* Z_a^\top Z_a) = \hat{\lambda}_s, \quad (39)$$

which yields

$$\text{tr}(Z_a^\top H_{\mathcal{I}^* \mathcal{I}^*}(X^*) Z_a) - \text{tr}(Z_a \Lambda^* Z_a^\top) \leq a^2 \hat{\lambda}_r + (1-a^2) \hat{\lambda}_s - \hat{\lambda}_s = a^2 (\hat{\lambda}_r - \hat{\lambda}_s). \quad (40)$$

The definition of  $Z_a$  gives

$$(Z_a^\top X_{\mathcal{I}^*}^*)_{pq} = \begin{cases} az^\top x_q, & \text{if } p = s, q \neq s; \\ \sqrt{1-a^2}, & \text{if } p = s, q = s; \\ 0, & \text{otherwise.} \end{cases} \quad (41)$$

Hence, we have

$$\begin{aligned} \text{tr}((X_{\mathcal{I}^*}^*)^\top Z_a \Lambda^* Z_a^\top X_{\mathcal{I}^*}^*) &= \text{tr}(\Lambda^* Z_a^\top X_{\mathcal{I}^*}^* (X_{\mathcal{I}^*}^*)^\top Z_a) = \hat{\lambda}_s \left( \sum_{q=1, q \neq s}^p a^2 (z^\top x_q)^2 + (1-a^2) \right) \\ &= \hat{\lambda}_s \left( \sum_{q=1}^p a^2 (z^\top x_q)^2 + (1-a^2) - a^2 (z^\top x_s)^2 \right) \\ &= \hat{\lambda}_s (1 + a^2 \|z^\top X_{\mathcal{I}^*}^*\|_2^2 - a^2) = a^2 \hat{\lambda}_s \|z^\top X_{\mathcal{I}^*}^*\|_2^2 + (1-a^2) \hat{\lambda}_s. \end{aligned} \quad (42)$$

and

$$\begin{aligned} \text{tr}(Z_a^\top X_{\mathcal{I}^*}^* \Lambda^* (X_{\mathcal{I}^*}^*)^\top Z_a) &= \left( \sum_{q=1, q \neq s}^p a^2 (z^\top x_q)^2 \hat{\lambda}_q + (1-a^2) \hat{\lambda}_s \right) \\ &\geq \left( \sum_{q=1, q \neq s}^p a^2 (z^\top x_q)^2 \hat{\lambda}_1 + (1-a^2) \hat{\lambda}_s \right) \\ &= \left( \sum_{q=1}^p a^2 (z^\top x_q)^2 \hat{\lambda}_1 + (1-a^2) \hat{\lambda}_s \right) \\ &= a^2 \hat{\lambda}_1 \|z^\top X_{\mathcal{I}^*}^*\|_2^2 + (1-a^2) \hat{\lambda}_s. \end{aligned} \quad (43)$$

Combining (42) and (43) together yields

$$\begin{aligned}
& \text{tr}((X_{\mathcal{I}^*}^*)^\top Z_a \Lambda^* Z_a^\top X_{\mathcal{I}^*}^*) - \text{tr}(Z_a^\top X_{\mathcal{I}^*}^* \Lambda^* (X_{\mathcal{I}^*}^*)^\top Z_a) \\
& \leq (a^2 \hat{\lambda}_s \|z^\top X_{\mathcal{I}^*}^*\|_2^2 + (1 - a^2) \hat{\lambda}_s) - (a^2 \hat{\lambda}_1 \|z^\top X_{\mathcal{I}^*}^*\|_2^2 + (1 - a^2) \hat{\lambda}_s) = a^2 (\hat{\lambda}_s - \hat{\lambda}_1) \|z^\top X_{\mathcal{I}^*}^*\|_2^2 \\
& \leq a^2 (\hat{\lambda}_s - \hat{\lambda}_1).
\end{aligned} \tag{44}$$

The equality (37) gives

$$\text{diag}(X_{\mathcal{I}^*}^* Z_a^\top \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp) = a \text{diag}(x_s z^\top) = \frac{a x_{is} x_{js}}{\sqrt{x_{is}^2 + x_{js}^2}} e_{(j,-i)}. \tag{45}$$

Let  $Z_a$  with  $a = \sqrt{\frac{\kappa}{\hat{\lambda}_r - \hat{\lambda}_1}}$ . Using (40) and (44), we have

$$\begin{aligned}
& \text{tr}(Z_a^\top H_{\mathcal{I}^* \mathcal{I}^*}(X_{\mathcal{I}^*}^*) Z_a) + \text{tr}((X_{\mathcal{I}^*}^*)^\top Z_a \Lambda_{\mathcal{I}^*} Z_a^\top X_{\mathcal{I}^*}^*) \\
& - \text{tr}(Z_a^\top X_{\mathcal{I}^*}^* \Lambda_{\mathcal{I}^*} (X_{\mathcal{I}^*}^*)^\top Z_a) - \text{tr}(Z_a \Lambda_{\mathcal{I}^*} Z_a^\top) \leq a^2 (\hat{\lambda}_r - \hat{\lambda}_1) = \kappa.
\end{aligned} \tag{46}$$

It follows from our definition of strong local minimizers that

$$\begin{aligned}
& \text{tr}(Z_a^\top H_{\mathcal{I}^* \mathcal{I}^*}(X^*) Z_a) + \text{tr}((X_{\mathcal{I}^*}^*)^\top Z_a \Lambda^* Z_a^\top X_{\mathcal{I}^*}^*) - \text{tr}(Z_a^\top X_{\mathcal{I}^*}^* \Lambda^* (X_{\mathcal{I}^*}^*)^\top Z_a) - \text{tr}(Z_a \Lambda^* Z_a^\top) \\
& + 2 \text{diag}(X_{\mathcal{I}^*}^* Z_a^\top \mathbf{P}_X^\perp)^\top \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} \text{Diag} \left( (\rho_{\mathcal{I}^*}^*)^{-\frac{2}{3}} \right) \right) \text{diag}(X_{\mathcal{I}^*}^* Z_a^\top \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp) \geq \kappa,
\end{aligned} \tag{47}$$

which together with (46) gives

$$\text{diag}(X_{\mathcal{I}^*}^* Z_a^\top \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp)^\top \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} \text{Diag} \left( (\rho_{\mathcal{I}^*}^*)^{-\frac{2}{3}} \right) \right) \text{diag}(X_{\mathcal{I}^*}^* Z_a^\top \mathbf{P}_{X_{\mathcal{I}^*}^*}^\perp) \geq 0. \tag{48}$$

Substituting (45) into (48), we obtain

$$e_{(j,-i)}^\top \left( L_{\mathcal{I}^* \mathcal{I}^*}^\dagger - \frac{\gamma}{3} \text{Diag} \left( (\rho_{\mathcal{I}^*}^*)^{-\frac{2}{3}} \right) \right) e_{(j,-i)} \geq 0. \tag{49}$$

Expanding the terms of (49) yields

$$(L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{jj} - 2(L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{ji} + (L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{ii} - \frac{\gamma}{3} (\rho_{\mathcal{I}^*}^*)_j^{-\frac{2}{3}} - \frac{\gamma}{3} (\rho_{\mathcal{I}^*}^*)_i^{-\frac{2}{3}} \geq 0, \tag{50}$$

which implies

$$(L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{jj} - 2(L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{ji} + (L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{ii} \geq \frac{\gamma}{3} (\rho_{\mathcal{I}^*}^*)_j^{-\frac{2}{3}}. \tag{51}$$

Therefore, we obtain

$$(\rho_{\mathcal{I}^*}^*)_j \geq \left( \frac{\gamma}{3((L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{jj} - 2(L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{ji} + (L_{\mathcal{I}^* \mathcal{I}^*}^\dagger)_{ii})} \right)^{\frac{3}{2}} \geq c_j, \tag{52}$$

where  $c_j$  is defined in (34). Similarly, we can prove (52) holds for any  $j \in \mathcal{J}$ . This completes the proof.  $\square$

### 3 Analysis of the KS Equation

#### 3.1 Formulating the KS Equation as a Fixed Point Map

The KS equation (5) constitutes a nonlinear system with respect to  $X$ . Note that the Hamiltonian matrix (4) is a symmetric matrix function with respect to  $\rho$  as

$$\hat{H}(\rho) := \frac{1}{2}L + V_{ion} + \text{Diag}(L^\dagger \rho) + \text{Diag}(\mu_{xc}(\rho)^T e), \quad (53)$$

and the KS equation becomes

$$\begin{cases} \hat{H}(\rho)X = X\Lambda, \\ X^T X = I, \end{cases} \quad (54)$$

where  $X \in \mathbb{R}^{n \times p}$  and  $\Lambda \in \mathbb{R}^{p \times p}$  is a diagonal matrix consisting of the  $p$  smallest eigenvalues of  $\hat{H}(\rho)$ . The eigenvalue decomposition of  $\hat{H}(\rho)$  is determined once  $\rho$  is given. Hence, we can write  $X$  as  $X(\rho)$  to reflect the dependence on  $\rho$  and the KS equation (5) can be viewed as a system of nonlinear equations with respect to the charge density  $\rho$  as

$$\rho = \text{diag}(X(\rho)X(\rho)^T). \quad (55)$$

Alternatively, the function

$$V := \mathcal{V}(\rho) = L^\dagger \rho + \mu_{xc}(\rho)^T e \quad (56)$$

is called potential and the Hamiltonian matrix  $\hat{H}(\rho)$ , by convenient abuse of notation, can be expressed as

$$H(V) := \frac{1}{2}L + V_{ion} + \text{Diag}(V). \quad (57)$$

Obviously, it holds  $\hat{H}(\rho) = H(\mathcal{V}(\rho))$ . Therefore,  $X$  can be interpreted as an implicit function of  $V$ . Let  $X(V) \in \mathbb{R}^{n \times p}$  be the eigenvectors corresponding to the  $p$  smallest eigenvalues of  $H(V)$ . Then, the fixed point map (55) is a system of nonlinear equations with respect to  $V$  as

$$\begin{cases} V = \mathcal{V}(F_\phi(V)), \\ F_\phi(V) = \text{diag}(X(V)X(V)^T). \end{cases} \quad (58)$$

The fixed point map (58) is well defined if there is a gap between the  $p$ th and  $(p+1)$ st smallest eigenvalues of  $H(V)$ . However, when these two eigenvalues are equal, there exists ambiguity on choosing the eigenvectors  $X(V)$  since the multiplicity is greater than one. A common approach is to revise  $F_\phi(V)$  in (58) by constructing a proper filter function. Let  $q_1(V), \dots, q_n(V)$  be the eigenvectors of  $H(V)$  associated with eigenvalues  $\lambda_1(V), \dots, \lambda_n(V)$ , respectively. A particular choice of the filter function is the Fermi-Dirac distribution of the form

$$f_\mu(t) := \frac{1}{1 + e^{\beta(t-\mu)}}, \quad (59)$$

where  $\mu$  is the solution of the equations

$$\sum_{i=1}^n f_\mu(\lambda_i(V)) = p. \quad (60)$$

Since the left hand side of (60) is monotonic with respect to  $\mu$  for a fixed  $\beta$ , the solution to (60) is unique for any

choice of  $\beta$  and  $\lambda_i$ . Then the fixed map (58) is replaced by the approximation

$$\begin{cases} V = \mathcal{V}(F_{f_\mu}(V)), \\ F_{f_\mu}(V) = \text{diag} \left( \sum_{i=1}^n f_\mu(\lambda_i(V)) q_i(V) q_i(V)^\text{T} \right). \end{cases} \quad (61)$$

### 3.2 The Jacobian of the Fixed Point Maps

We first reformulate the functions  $F_\phi(V)$  in (58) and  $F_{f_\mu}(V)$  in (61) as the form of spectral operators. Using the differentiability of spectral operators, they can be proved to be differentiable under some conditions. Let  $\{\lambda_i(V), q_i(V)\}$  be the eigenpairs of  $H(V)$  and assume that the eigenvalues  $\lambda_1(V), \dots, \lambda_n(V)$  are sorted in an ascending order,

$$\lambda_1(V) \leq \dots \leq \lambda_p(V) \leq \lambda_{p+1}(V) \leq \dots \leq \lambda_n(V).$$

The eigenvalue decomposition of  $H(V)$  can be written as

$$H(V) = Q(V)\Pi(V)Q(V)^\text{T}, \quad (62)$$

where  $Q(V)$  and  $\Pi(V)$  are

$$Q(V) = [q_1(V), q_2(V), \dots, q_n(V)] \in \mathbb{R}^{n \times n} \quad \text{and} \quad \Pi(V) = \text{Diag}(\lambda_1(V), \lambda_2(V), \dots, \lambda_n(V)) \in \mathbb{R}^{n \times n}. \quad (63)$$

Hence, the function  $F_\phi(V)$  in (58) is equivalent to

$$F_\phi(V) = \text{diag}(Q(V)\phi(\Pi(V))Q(V)^\text{T}), \quad (64)$$

where  $\phi(\Pi) = \text{Diag}(\phi(\lambda_1(V)), \phi(\lambda_2(V)), \dots, \phi(\lambda_n(V)))$  and

$$\phi(t) := \begin{cases} 1 & \text{for } t \leq \frac{\lambda_p(V) + \lambda_{p+1}(V)}{2}, \\ 0 & \text{for } t > \frac{\lambda_p(V) + \lambda_{p+1}(V)}{2}. \end{cases} \quad (65)$$

Similarly, the function  $F_{f_\mu}(V)$  in (61) in the spectral operator form is

$$F_{f_\mu}(V) = \text{diag}(Q(V)f_\mu(\Pi(V))Q(V)^\text{T}). \quad (66)$$

Let  $\mu_1, \dots, \mu_{r(V)}$  be the distinct eigenvalues among  $\{\lambda_1(V), \dots, \lambda_n(V)\}$ ,  $r(V)$  be the total number of distinct values and  $r_p(V)$  be the number of distinct eigenvalues no greater than  $\lambda_p$ . For any  $k = 1, \dots, r(V)$ , the set of indices  $i$  such that  $\lambda_i = \mu_k$  is denoted by  $\alpha_k := \{i \mid \lambda_i = \mu_k, i = 1, \dots, n\}$ . The next lemma shows the directional derivative of  $F_\phi(V)$  by using the differentiability of the spectral operators [11, 14, 22, 28, 27].

**Lemma 3.1.** *Suppose that Assumption 2.5 holds at  $H(V)$ , i.e.,  $\lambda_{p+1}(V) > \lambda_p(V)$ . Then  $F_\phi(V)$  is continuously differentiable and its directional derivative at  $V$  along  $z \in \mathbb{R}^n$  is*

$$\partial_V F_\phi(V)[z] = \text{diag} \left( Q(V) (g_\phi(\Pi(V)) \circ (Q(V)^\text{T} \text{Diag}(z) Q(V))) Q(V)^\text{T} \right), \quad (67)$$

where “ $\circ$ ” denotes the Hadamard product between two matrices, and  $g_\phi(\Pi(V)) \in \mathbb{R}^{n \times n}$  is the so-called first divided

difference matrix defined as

$$(g_\phi(\Pi(V)))_{ij} = \begin{cases} \frac{1}{\lambda_i(V) - \lambda_j(V)} & \text{if } i \in \alpha_k, j \in \alpha_l, k \leq r_p(V), l > r_p(V), \\ \frac{-1}{\lambda_i(V) - \lambda_j(V)} & \text{if } i \in \alpha_k, j \in \alpha_l, k > r_p(V), l \leq r_p(V), \\ 0 & \text{otherwise.} \end{cases} \quad (68)$$

*Proof.* The chain rule gives

$$\partial_V F_\phi(V)[z] = \frac{d \text{diag}(Q\phi(\Pi)Q^T)}{dH} [\partial_V H(V)[z]]. \quad (69)$$

By applying the continuous differentiability of the spectral operators in Proposition 2.10 of [14], the function  $Q\phi(\Pi)Q^T$  is differentiable with respect to  $H$  and its directional derivative is given by

$$\frac{dQ\phi(\Pi)Q^T}{dH}[S] = Q(g_\phi(\Pi) \circ (Q^T S Q)) Q^T, \quad \text{for all } S \in \mathbb{S}^n, \quad (70)$$

where, for any  $i, j = 1, \dots, n$ ,

$$(g_\phi(\Pi(V)))_{ij} = \begin{cases} \frac{\phi(\lambda_i(V)) - \phi(\lambda_j(V))}{\lambda_i(V) - \lambda_j(V)} & \text{if } i \in \alpha_k, j \in \alpha_l, k \neq l, \\ 0 & \text{otherwise.} \end{cases} \quad (71)$$

Substituting (65) into (71) yields the specific form of  $g_\phi(\pi(V))$  in (68). Since  $\text{diag}(\cdot)$  is a linear function, we obtain

$$\begin{aligned} \frac{d \text{diag}(Q\phi(\Lambda)Q^T)}{dH}[S] &= \frac{d \text{diag}(Q\phi(\Lambda)Q^T)}{dQ\phi(\Lambda)Q^T} \frac{dQ\phi(\Lambda)Q^T}{dH}[S] \\ &= \text{diag}(Q(g_\phi(\Pi) \circ (Q^T S Q)) Q^T), \quad \text{for all } S \in \mathbb{S}^n. \end{aligned} \quad (72)$$

It follows from (57) that

$$\partial_V H(V)[z] = \text{Diag}(z). \quad (73)$$

Plugging (72) and (73) into (69), we obtain (67). This completes the proof.  $\square$

**Remark 3.2.** Computing  $\partial_V F_\phi(V)[z]$  requires all the eigenvectors  $Q(V)$  and all eigenvalues  $\Pi(V)$ . Let  $E_{j,p}$  ( $O_{j,p}$ ) be the  $j \times p$  matrix with ones (zeros) at all its entries. Then the matrix  $g_\phi(\Pi(V)) \in \mathbb{R}^{n \times n}$  takes the specific form

$$g_\phi(\Pi(V)) = \begin{pmatrix} O_{p,p} & G \\ G^T & O_{n-p, n-p} \end{pmatrix},$$

where

$$G = \begin{pmatrix} \frac{1}{\mu_1 - \mu_{r_p(V)+1}} E_{|\alpha_1|, |\alpha_{r_p(V)+1}|} & \cdots & \frac{1}{\mu_1 - \mu_r(V)} E_{|\alpha_1|, |\alpha_r(V)|} \\ \vdots & \ddots & \vdots \\ \frac{1}{\mu_{r_p(V)} - \mu_{r_p(V)+1}} E_{|\alpha_{r_p(V)}|, |\alpha_{r_p(V)+1}|} & \cdots & \frac{1}{\mu_{r_p(V)} - \mu_r(V)} E_{|\alpha_{r_p(V)}|, |\alpha_r(V)|} \end{pmatrix}.$$

The directional derivative of  $F_{f_\mu}(V)[z]$  can be assembled in a similar fashion.

**Lemma 3.3.** *The function  $F_{f_\mu}(V)$  is continuously differentiable and its directional derivative at  $V$  along  $z \in \mathbb{R}^n$  is*

$$\partial_V F_{f_\mu}(V)[z] = \text{diag} \left( Q(V) \left( g_{f_\mu}(\Pi(V)) \circ \left( Q(V)^T \text{Diag}(z) Q(V) \right) \right) Q(V)^T \right), \quad (74)$$

where  $g_{f_\mu}(\Pi(V)) \in \mathbb{R}^{n \times n}$  is defined as, for any  $i, j = 1, \dots, n$ ,

$$(g_{f_\mu}(\Pi(V)))_{ij} = \begin{cases} \frac{f_\mu(\lambda_i(V)) - f_\mu(\lambda_j(V))}{\lambda_i(V) - \lambda_j(V)} & \text{if } i \in \alpha_k, j \in \alpha_l, k \neq l, \\ f'_\mu(\lambda_i(V)) & \text{otherwise.} \end{cases} \quad (75)$$

We next compute the Jacobian of  $\mathcal{V}(F_\phi(V))$  and  $\mathcal{V}(F_{f_\mu}(V))$ .

**Theorem 3.4.** *Let  $J(\rho)$  be defined as (8).*

1. *Suppose that Assumption 2.5 holds at  $H(V)$ , i.e.,  $\lambda_{p+1}(V) > \lambda_p(V)$ . Then the Jacobian of  $\mathcal{V}(F_\phi(V))$  at  $V$  is*

$$\partial_V \mathcal{V}(F_\phi(V))[z] = J(F_\phi(V)) \partial_V F_\phi(V)[z], \quad \text{for all } z \in \mathbb{R}^n. \quad (76)$$

2. *The Jacobian of  $\mathcal{V}(F_{f_\mu}(V))$  at  $V$  is*

$$\partial_V \mathcal{V}(F_{f_\mu}(V))[z] = J(F_{f_\mu}(V)) \partial_V F_{f_\mu}(V)[z], \quad \text{for all } z \in \mathbb{R}^n. \quad (77)$$

*Proof.* Note that

$$\partial_\rho(\mathcal{V}(\rho))[z] = J(\rho)z, \quad \text{for all } z \in \mathbb{R}^n. \quad (78)$$

Applying the chain rules to  $\partial_V \mathcal{V}(F_\phi(V))[z]$  and using (78) and (67), we obtain (76). This completes the proof.  $\square$

## 4 Convergence of the SCF iteration

### 4.1 The SCF Iteration and the Simple Mixing Scheme

Starting from an initial vector  $V^0 \in \mathbb{R}^n$ , the SCF iteration for solving the fixed point map (58) recursively computes the eigenpairs  $\{X(V^{i+1}), \Lambda(V^{i+1})\}$  as the solution of the linear eigenvalue problem:

$$\begin{aligned} H(V^i)X(V^{i+1}) &= X(V^{i+1})\Lambda(V^{i+1}), \\ X(V^{i+1})^T X(V^{i+1}) &= I, \end{aligned}$$

and then the potential is updated as

$$V^{i+1} = \mathcal{V}(F_\phi(V^i)). \quad (79)$$

When the difference between  $V^i$  and  $V^{i+1}$  is negligible, the system is said to be self-consistent and the SCF iteration is terminated.

The SCF iteration often converges slowly or even fails to converge. One of the heuristics for accelerating and stabilizing the SCF iteration is charge or potential mixing [17, 19]. Basically, the new potential  $V^{i+1}$  is constructed from a linear combination of the previously computed potential and the one obtained from certain schemes at current

iteration. In particular, the simple mixing scheme replaces (79) by updating

$$V^{i+1} = V^i - \alpha(V^i - \mathcal{V}(F_\phi(V^i))), \quad (80)$$

where  $\alpha$  is a properly chosen step size. Similarly, the SCF iteration using simple mixing for solving the fixed point map (61) is

$$V^{i+1} = V^i - \alpha(V^i - \mathcal{V}(F_{f_\mu}(V^i))). \quad (81)$$

## 4.2 Global Convergence Analysis

We first make the following assumptions.

**Assumption 4.1.** *The second-order derivatives of the exchange correlation functional  $\epsilon_{xc}(\rho)$  is uniformly bounded from above. Without loss of generality, we assume that there exists a constant  $\theta$  such that*

$$\|\partial\mu_{xc}(\rho)e\|_2 \leq \theta, \quad \text{for all } \rho \in \mathbb{R}^n. \quad (82)$$

Although we cannot verify Assumption 4.1 for any  $X \in \mathbb{R}^{n \times p}$ , it holds at a strong local minimizer using our lower bounds for nonzero charge densities in subsection 2.5 if the exchange correlation energy is (16).

It can be verified from the definition of the operator  $\partial_V F_\phi(V)[\cdot]$  in (67) that it is a linear map. The induced  $\ell_2$ -norm of  $\partial_V \mathcal{V}(F_\phi(V))$  and  $\partial_V F_\phi(V)[\cdot]$  are defined as

$$\|\partial_V \mathcal{V}(F_\phi(V))\|_2 = \max_{z \neq 0} \frac{\|\partial_V \mathcal{V}(F_\phi(V))[z]\|_2}{\|z\|_2} \quad \text{and} \quad \|\partial_V F_\phi(V)\|_2 = \max_{z \neq 0} \frac{\|\partial_V F_\phi(V)[z]\|_2}{\|z\|_2}, \quad (83)$$

respectively. The next lemma shows that their  $\ell_2$ -norms are bounded if Assumption 2.5 holds at  $H(V)$ .

**Lemma 4.2.** *If Assumption 2.5 holds at  $H(V)$  for a given  $V \in \mathbb{R}^n$ , then it holds*

$$\|\partial_V F_\phi(V)\|_2 \leq \frac{1}{\delta} \quad \text{and} \quad \|\partial_V \mathcal{V}(F_\phi(V))\|_2 \leq \frac{\|L^\dagger\|_2 + \theta}{\delta}. \quad (84)$$

*Proof.* For any  $z \in \mathbb{R}^n$ , we obtain

$$\begin{aligned} \|\partial_V F_\phi(V)[z]\|_2 &= \|\text{diag}(Q(V)(g_\phi(\Pi(V)) \circ (Q(V)^T \text{Diag}(z) Q(V))) Q(V)^T)\|_2 \\ &\leq \|Q(\rho)(g_\phi(\Pi(\rho)) \circ (Q(\rho)^T \text{Diag}(z) Q(\rho))) Q(\rho)^T\|_F \\ &= \|g_\phi(\Pi(\rho)) \circ (Q(\rho)^T \text{Diag}(z) Q(\rho))\|_F \\ &\leq \frac{1}{\delta} \|Q(\rho)^T \text{Diag}(z) Q(\rho)\|_F \\ &\leq \frac{1}{\delta} \|z\|_2, \end{aligned} \quad (85)$$

where the second inequality is due to  $|(g_\phi(\Pi(\rho)))_{ij}| \leq 1/\delta$ . Then the first inequality in (84) holds from the definitions (83) and (85). It follows from (76) and (85) that

$$\|\partial_V \mathcal{V}(F_\phi(V))[z]\|_2 \leq \|J(F_\phi(V))\|_2 \|\partial_V F_\phi(V)[z]\|_2 \leq \frac{\|L^\dagger\|_2 + \theta}{\delta} \|z\|_2. \quad (86)$$

This completes the proof.  $\square$

The set  $\{H(V) \mid V \in \mathbb{R}^n\}$  is called uniformly well posed (UWP) [2, 30] with respect to a constant  $\delta > 0$  if Assumption 2.5 holds at  $H(V)$  with  $\delta$  for any  $V \in \mathbb{R}^n$ . We next establish the convergence of the simple mixing scheme (80) when UWP holds.

**Theorem 4.3.** *Suppose that Assumption 4.1 holds and  $\{H(V) \mid V \in \mathbb{R}^n\}$  is UWP with a constant  $\delta$  such that*

$$b_1 := 1 - \frac{\|L^\dagger\|_2 + \theta}{\delta} > 0. \quad (87)$$

Let  $\{V^i\}$  be a sequence generated by the simple mixing scheme (80) using a step size  $\alpha$  satisfying

$$0 < \alpha < \frac{2}{2 - b_1}. \quad (88)$$

Then  $\{V^i\}$  converges to a solution of the KS equation (5) with linear convergence rate no more than  $|1 - \alpha| + \alpha(1 - b_1)$ .

*Proof.* For any  $V^i$ , it follows from (86), (87) and (88) that

$$\begin{aligned} & \|(1 - \alpha)I + \alpha \partial_V \mathcal{V}(F_\phi(V^i))\|_2 \\ & \leq |1 - \alpha| + |\alpha| \|\partial_V \mathcal{V}(F_\phi(V^i))\|_2 \\ & \leq \begin{cases} 1 - \alpha + \alpha \frac{\|L^\dagger\|_2 + \theta}{\delta} = 1 - \alpha b_1, & \text{if } 0 < \alpha < 1 \\ \alpha - 1 + \alpha \frac{\|L^\dagger\|_2 + \theta}{\delta} = \alpha(2 - b_1) - 1, & \text{if } \alpha \geq 1 \end{cases} \\ & < 1, \end{aligned}$$

which completes the proof.  $\square$

**Remark 4.4.** *When the step size  $\alpha = 1$ , the simple mixing scheme (80) becomes the SCF iteration (79) with the convergence rate  $\frac{\|L^\dagger\|_2 + \theta}{\delta}$ . Since neither  $p$  nor  $n$  is involved in (87), it is much weaker than  $\frac{12k\sqrt{n}\|L^\dagger\|_2 + \theta}{\delta} < 1$  required by Theorem 1 in [24].*

We next establish convergence to the solutions of the modified fixed-point map (61) without assuming the UWP properties.

**Theorem 4.5.** *Suppose that Assumption 4.1 holds and*

$$b_2 := 1 - \frac{\beta(\|L^\dagger\|_2 + \theta)}{4} > 0. \quad (89)$$

Let  $\{V^i\}$  be a sequence generated by the simple mixing scheme (81) using a step size  $\alpha$  satisfying

$$0 < \alpha < \frac{2}{2 - b_2}. \quad (90)$$

Then the sequence  $\{V^i\}$  converges to a solution of (61) with linear convergence rate no less than  $|1 - \alpha| + \alpha(1 - b_2)$ .

*Proof.* Using the mean value theorem and the fact that

$$|f'_\mu(t)| = \left| \frac{-\beta e^{\beta(t-\mu)}}{(1 + e^{\beta(t-\mu)})^2} \right| \leq \frac{\beta}{4},$$

we obtain  $|(g_{f_\mu}(\Pi(V)))_{ij}| \leq \beta/4$ , which yields

$$\|\partial_V \mathcal{V}(F_{f_\mu}(V))\|_2 \leq \frac{\beta(\|L^\dagger\|_2 + \theta)}{4}.$$

Then, the convergence of (81) is proved similar to that of Theorem 4.3.  $\square$

**Remark 4.6.** Suppose that UWP holds and  $f_\mu$  is chosen such that

$$\begin{cases} \frac{1}{1+e^{\beta(\lambda_p-\mu)}} \geq 1-\gamma, \\ \frac{1}{1+e^{\beta(\lambda_{p+1}-\mu)}} \leq \gamma, \end{cases} \quad (91)$$

where  $\gamma \ll 1$  is a constant. It can be shown that  $\beta \geq \frac{2}{\delta} \cdot \ln \frac{1-\gamma}{\gamma}$ . Hence, we have  $\frac{\beta}{4} \geq \frac{1}{8}$  and the condition (87) is implied by (89) when  $\ln \frac{1-\gamma}{\gamma} \geq 2$  or equivalently  $\gamma \leq \frac{1}{e^2+1} \approx 0.12$ . On the other hand, the closer  $\gamma$  is to zero, the closer  $f_\mu$  is to  $\phi$  from (91). Therefore, the convergence rate of the fixed-point iteration using  $F_\phi$  is better than that of  $F_{f_\mu}$  when  $F_{f_\mu}$  is sufficiently close to  $F_\phi$ .

**Remark 4.7.** The convergence of the SCF iteration without simple mixing for solving a special KS equation without the exchange correlation energy is established in [30] under the condition

$$\frac{n^4 \beta \|L^\dagger\|_2}{2} < 1. \quad (92)$$

We can see that our condition is weaker than (92) since  $n^4$  is not required.

### 4.3 Local Convergence Analysis

Suppose that  $V^*$  is a solution of the fixed point map (80). Let  $B(V^*, \eta) := \{V \mid \|V - V^*\|_2 \leq \eta\}$  be a neighborhood of  $V^*$  for a given  $\eta > 0$ . The Taylor expansion at  $V^*$  yields

$$\begin{aligned} V^{k+1} - V^* &= V^k - \alpha(V^k - \mathcal{V}(F_\phi(V^k))) - (V^* - \alpha(V^* - \mathcal{V}(F_\phi(V^*)))) \\ &= (I - \alpha(I - \partial_V \mathcal{V}(F_\phi(V^*))))[V^k - V^*] + o(\|V^k - V^*\|_2), \quad \text{for all } V^k \in B(V^*, \eta). \end{aligned} \quad (93)$$

If the spectral radius of the operator  $I - \alpha(I - \partial_V \mathcal{V}(F_\phi(V^*)))$  is less than one, there must exist a sufficiently small  $\eta$  so that the simple mixing scheme (80) initiating from a point in  $B(V^*, \eta)$  converges to  $V^*$  linearly.

We first present a few properties of the linear operators. Denote the space of linear operators by

$$\mathbb{L}(\mathbb{R}^n, \mathbb{R}^n) := \{\mathcal{P} \mid \mathcal{P} : \mathbb{R}^n \mapsto \mathbb{R}^n \text{ is a linear map}\}.$$

Since  $\mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$  is isomorphic to  $\mathbb{R}^{n \times n}$ , the eigenvalue, eigenvector and the spectrum for any linear operator can be defined similar to a matrix. For a given  $\mathcal{P} \in \mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$ , if a scalar  $\lambda \in \mathbb{C}$  and a nonzero vector  $z \in \mathbb{C}^n$  satisfy

$$\mathcal{P}[z] = \lambda z, \quad (94)$$

the scalar  $\lambda$  and the vector  $z$  are called the eigenvalue and eigenvector of  $\mathcal{P}$ , respectively. The spectrum of  $\mathcal{P}$ , denoted by  $\lambda(\mathcal{P})$ , is the set consisting of all the eigenvalues of  $\mathcal{P}$ . The spectral radius, denoted by  $\varrho(\mathcal{P})$ , is the largest absolute value of all elements in its spectrum. The operator  $\mathcal{P}$  is called symmetric if  $y^T \mathcal{P}[x] = x^T \mathcal{P}[y]$  for any  $x, y \in \mathbb{R}^n$ .

**Definition 4.8.** Given  $\mathcal{P} \in \mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$ , the matrix  $P = (\mathcal{P}[e_1], \dots, \mathcal{P}[e_n])$  is called the basic transformation matrix of  $\mathcal{P}$ , where  $e_i, i = 1, 2, \dots, n$ , is the  $i$ th column of the identity matrix. A linear operator  $\mathcal{P}^* \in \mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$  is called the adjoint operator of  $\mathcal{P}$  if  $\mathcal{P}^*[x] = P^T x$  holds for all  $x \in \mathbb{R}^n$ .

Let  $P$  be the basic transformation matrix of  $\mathcal{P} \in \mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$ . Then  $P$  is symmetric if and only if  $\mathcal{P}$  is symmetric. Moreover,  $\mathcal{P}$  and  $P$  has the same spectrum since  $\mathcal{P}[z] = Pz$ . Let  $M_1, M_2 \in \mathbb{R}^{n \times n}$  be two real matrices and  $P_1$  and  $P_2$  be the basic transformation matrices of  $\mathcal{P}_1, \mathcal{P}_2 \in \mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$ , respectively. Then  $M_1 P_1 + M_2 P_2$  is the basic transformation matrix of the linear operator  $M_1 \mathcal{P}_1 + M_2 \mathcal{P}_2$ . A linear operator  $\mathcal{P} \in \mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$  is called positive semidefinite if  $z^T \mathcal{P}[z] \geq 0$  for all  $z \in \mathbb{R}^n$ . We next show that the eigenvalues of the product of a symmetric matrix and a symmetric positive semidefinite linear operator are real.

**Lemma 4.9.** Suppose that  $M \in \mathbb{R}^{n \times n}$  is a symmetric matrix, and  $\mathcal{P} \in \mathbb{L}(\mathbb{R}^n, \mathbb{R}^n)$  is a symmetric positive semidefinite linear operator. Then all the eigenvalues of the linear operator  $M\mathcal{P}$  are real. Furthermore, it holds

$$\lambda_{\max}(M\mathcal{P}) \leq \begin{cases} \lambda_{\max}(M)\lambda_{\max}(\mathcal{P}), & \text{if } \lambda_{\max}(M) \geq 0, \\ \lambda_{\max}(M)\lambda_{\min}(\mathcal{P}), & \text{otherwise,} \end{cases} \quad (95)$$

$$\lambda_{\min}(M\mathcal{P}) \geq \begin{cases} \lambda_{\min}(M)\lambda_{\min}(\mathcal{P}), & \text{if } \lambda_{\min}(M) \geq 0, \\ \lambda_{\min}(M)\lambda_{\max}(\mathcal{P}), & \text{otherwise.} \end{cases} \quad (96)$$

*Proof.* Let  $P$  be the basic transformation matrix of  $\mathcal{P}$ . It suffices to prove the statements with  $\mathcal{P}$  replaced by  $P$ . Since  $\mathcal{P}$  is symmetric positive semidefinite,  $P$  is also symmetric positive semidefinite. Hence, it can be diagonalized as  $P = UDU^T$ , where  $U$  is orthogonal and  $D = \text{Diag}(\mu_1, \dots, \mu_n)$  such that  $\mu_i \geq 0$ . Define  $D^{\frac{1}{2}} := \text{Diag}(\mu_1^{\frac{1}{2}}, \dots, \mu_n^{\frac{1}{2}})$  and write  $P^{\frac{1}{2}} = UD^{\frac{1}{2}}U^T$ . Then we obtain  $P = P^{\frac{1}{2}}P^{\frac{1}{2}}$ . We now prove that every eigenvalue of  $R := P^{\frac{1}{2}}MP^{\frac{1}{2}}$  is an eigenvalue of  $MP$  and vice versa. It is known that the eigenvalues of a matrix are continuous functions of the matrix entries. Let  $D_\epsilon := D + \epsilon I$  and  $P_\epsilon := UD_\epsilon U^T$  for  $\epsilon \geq 0$ . Then  $P_\epsilon \rightarrow P$  and  $P_\epsilon^{\frac{1}{2}} := UD_\epsilon^{\frac{1}{2}}U^T \rightarrow P^{\frac{1}{2}}$  as  $\epsilon \rightarrow 0$ . Hence,  $MP_\epsilon \rightarrow MP$  and  $R_\epsilon := P_\epsilon^{\frac{1}{2}}MP_\epsilon^{\frac{1}{2}} \rightarrow R$  as  $\epsilon \rightarrow 0$ . Since  $P_\epsilon^{\frac{1}{2}}$  is invertible, we have  $R_\epsilon = P_\epsilon^{\frac{1}{2}}MP_\epsilon P_\epsilon^{-\frac{1}{2}}$ . Therefore,  $R_\epsilon$  and  $MP_\epsilon$  have the same eigenvalues. As  $\epsilon \rightarrow 0$ , these eigenvalues converge to those of  $R$  and  $MP$ , respectively. Hence,  $R$  and  $MP$  have the same eigenvalues. The symmetry of  $R$  further implies that the eigenvalues of  $MP$  are real.

Since  $\lambda_{\max}(M)I \succeq M$ , we obtain

$$\lambda_{\max}(M)P = P^{\frac{1}{2}}(\lambda_{\max}(M) - M)P^{\frac{1}{2}} + P^{\frac{1}{2}}MP^{\frac{1}{2}} \succeq P^{\frac{1}{2}}MP^{\frac{1}{2}},$$

which yields (95) since the eigenvalues of  $R = P^{\frac{1}{2}}MP^{\frac{1}{2}}$  and  $MP$  are the same. Similarly, (96) holds due to  $M \succeq \lambda_{\min}(M)I$  and

$$P^{\frac{1}{2}}MP^{\frac{1}{2}} = P^{\frac{1}{2}}(M - \lambda_{\min}(M))P^{\frac{1}{2}} + \lambda_{\min}(M)P \succeq \lambda_{\min}(M)P.$$

This completes the proof.  $\square$

The next lemma shows that  $\partial_V F_\phi(V)[\cdot]$  is negative semidefinite.

**Lemma 4.10.** For any  $z \in \mathbb{R}^n$ , it holds  $z^T \partial_V F_\phi(V)[z] \leq 0$ .

*Proof.* For any  $z \in \mathbb{R}^n$ , we have

$$z^T \partial_V F_\phi(V)[z] = z^T \text{diag} (Q(V) (g_\phi(\Pi(V)) \circ (Q(V)^T \text{Diag}(z) Q(V))) Q(V)^T)$$

$$\begin{aligned}
&= \langle (Q(V)^T \text{Diag}(z) Q(V)), g_\phi(\Pi(V)) \circ (Q(V)^T \text{Diag}(z) Q(V)) \rangle \\
&= e^T (g_\phi(\Pi(V)) \circ (Q(V)^T \text{Diag}(z) Q(V)) \circ (Q(V)^T \text{Diag}(z) Q(V))) e \\
&\leq 0,
\end{aligned}$$

where the third equality uses the properties of the Hadamard products and the inequality is due to

$$(Q(V)^T \text{Diag}(z) Q(V)) \circ (Q(V)^T \text{Diag}(z) Q(V)) \geq 0 \text{ and } g_\phi(\Pi(V)) \leq 0.$$

This completes the proof.  $\square$

We now establish the local convergence result for the simple mixing scheme.

**Theorem 4.11.** *Let  $V^*$  be a solution of the KS equation (5). Suppose that Assumption 4.1 holds and Assumption 2.5 is valid at  $H(V^*)$  with a constant  $\delta$  satisfying*

$$\delta > -\lambda_{\min}^*, \quad (97)$$

where  $\lambda_{\min}^* := \min\{0, \lambda_{\min}(J(F_\phi(V^*)))\}$ . There exists an open neighborhood  $\Omega$  of  $V^*$ , such that the sequence  $\{V^i\}$  generated by the simple mixing scheme (80) using  $V^0 \in \Omega$  and a step size

$$\alpha \in \left(0, \frac{2\delta}{\|L^\dagger\|_2 + \theta + \delta}\right) \quad (98)$$

converges to  $V^*$  with  $R$ -linear convergence rate no more than

$$\max \left\{ \left(1 - \alpha \frac{\delta + \lambda_{\min}^*}{\delta}\right), \left(\alpha \frac{\|L^\dagger\|_2 + \theta + \delta}{2\delta} - 1\right) \right\}.$$

*Proof.* The Taylor expansion (93) implies that local convergence of the scheme (80) holds if

$$\varrho(I - \alpha \mathcal{A}) < 1, \quad (99)$$

where  $\mathcal{A} := I - J(F_\phi(V^*))\partial_V F_\phi(V^*)$ . According to Lemma 4.10,  $-\partial_V F_\phi(V^*)$  is symmetric positive semidefinite. Using Lemma 4.9, we conclude that all the eigenvalues of  $\mathcal{A}$  are real. Hence, (99) is guaranteed if

$$\lambda_{\min}(\mathcal{A}) > 0; \quad (100)$$

$$\alpha \lambda_{\max}(\mathcal{A}) < 2. \quad (101)$$

Note that  $\lambda_{\min}(\mathcal{A}) = 1 + \lambda_{\min}(J(F_\phi(V^*))(-\partial_V F_\phi(V^*)))$ . Using Lemma 4.9,  $\lambda_{\max}(-\partial_V F_\phi(V^*)) \leq \frac{1}{\delta}$  from Lemma 4.2 and the definition of  $\lambda_{\min}^*$ , we obtain

$$\begin{aligned}
\lambda_{\min}(\mathcal{A}) - 1 &\geq \begin{cases} \lambda_{\min}(J(F_\phi(V^*)))\lambda_{\min}(-\partial_V F_\phi(V^*)), & \text{if } \lambda_{\min}(J(F_\phi(V^*))) \geq 0, \\ \lambda_{\min}(J(F_\phi(V^*)))\lambda_{\max}(-\partial_V F_\phi(V^*)), & \text{otherwise} \end{cases} \\
&\geq \begin{cases} 0, & \text{if } \lambda_{\min}(J(F_\phi(V^*))) \geq 0, \\ \frac{1}{\delta}\lambda_{\min}(J(F_\phi(V^*))), & \text{otherwise} \end{cases} \\
&\geq \frac{\lambda_{\min}^*}{\delta},
\end{aligned}$$

which yields (100) from the assumption (97).

Using Lemma 4.9 again, we have

$$\begin{aligned}\lambda_{\max}(\mathcal{A}) &\leq 1 + \lambda_{\max}(J(F_\phi(V^*))(-\partial_V F_\phi(V^*))) \\ &\leq 1 + \max\{0, \lambda_{\max}(J(F_\phi(V^*)))\lambda_{\max}(-\partial_V F_\phi(V^*))\} \leq 1 + \frac{\|L^\dagger\|_2 + \theta}{\delta},\end{aligned}\quad (102)$$

which together with (98) gives (101).  $\square$

The condition (97) can be much weaker than  $\|L^\dagger\|_2 + \theta < \delta$  required in Theorem 4.3.

**Corollary 4.12.** *Suppose that Assumption 4.1 holds. Then the condition (97) holds if*

$$\max(\theta - \lambda_{\min}(L^\dagger), 0) < \delta. \quad (103)$$

*Proof.* It follows from (8) and Assumption 4.1 that

$$\lambda_{\min}(J(F_\phi(V))) = \lambda_{\min}(L^\dagger + \partial_{\mu_{xc}}(F_\phi(V))e) \geq \lambda_{\min}(L^\dagger) + \lambda_{\min}(\mu_{xc}(F_\phi(V))e) \geq \lambda_{\min}(L^\dagger) - \theta.$$

Hence, (97) holds from the definition of  $\lambda_{\min}^*$ .  $\square$

In particular, when  $J(F_\phi(V^*))$  is positive semidefinite, we have  $\lambda_{\min}^* = 0$  and (97) is a direct consequence of Assumption 2.5.

**Corollary 4.13.** *Suppose that Assumption 2.5 holds at  $H(V^*)$  and  $J(F_\phi(V^*))$  is positive semidefinite. Then the condition (97) holds.*

We can obtain the following local convergence result for the modified fixed-point map (61) in the same manner as Theorem 4.5.

**Corollary 4.14.** *Suppose that Assumption 4.1 holds and*

$$\frac{4}{\beta} > -\lambda_{\min}^*, \quad (104)$$

where  $\lambda_{\min}^* := \min\{0, \lambda_{\min}(J(F_\phi(V^*)))\}$ . Let  $V^*$  be a solution of the KS equation (5). There exists an open neighborhood  $\Omega$  of  $V^*$ , such that the sequence  $\{V^i\}$  generated by the simple mixing scheme (81) using  $V^0 \in \Omega$  and a step size

$$\alpha \in \left(0, \frac{8}{(\|L^\dagger\|_2 + \theta)\beta + 4}\right) \quad (105)$$

converges to  $V^*$  with  $R$ -linear convergence rate no more than

$$\max\left\{\left(1 - \alpha \frac{\lambda_{\min}^* \beta + 4}{4}\right), \left(\alpha \frac{(\|L^\dagger\|_2 + \theta)\beta + 4}{8} - 1\right)\right\}.$$

## 5 Convergence Analysis of Approximate Newton Approaches

The generalized Jacobian  $\partial_V \mathcal{V}(F(V))$  in (76) suggests that Newton's method for solving the fixed point map (58) is

$$V^{i+1} = V^i - \alpha (I - J(F_\phi(V^i))\partial_V F_\phi(V^i))^{-1} (V^i - \mathcal{V}(F_\phi(V^i))),$$

where  $\alpha$  is a step size. Obviously, this method is not computationally practical for solving the fixed-point maps due to the presence of all eigenvectors and eigenvalues in  $\partial_V F_\phi(V)[\cdot]$ . In this section, we propose two approximate Newton approaches in the form

$$V^{i+1} = V^i - \alpha (I - D^i)^{-1} (V^i - \mathcal{V}(F_\phi(V^i))), \quad (106)$$

where  $\alpha > 0$  and  $D^i \in \mathbb{R}^{n \times n}$  is a matrix for approximating the Jacobian  $\partial_V \mathcal{V}(F(V^i))$ .

**Theorem 5.1.** *Suppose that Assumption 4.1 and UWP hold. Let  $\{V^i\}$  be a sequence generated by (106) using  $\{D^i\}$  and a step size  $\alpha$  such that*

$$0 < \alpha < \frac{2}{b_2}, \quad 0 < \gamma_{\min} \leq \sigma_{\min}(I - D^i) \text{ and } \sigma_{\max}(I - D^i) \leq \gamma_{\max},$$

where  $b_2 := 1 + \frac{\|L^\dagger\|_2 + \theta}{\delta}$ , and  $\sigma_{\min}$  and  $\sigma_{\max}$  are the smallest and largest singular values of  $I - D^i$ , respectively. If  $b_1 := 1 - \frac{\gamma_{\max} \|L^\dagger\|_2 + \theta}{\gamma_{\min} \delta} > 0$ , then  $\{V^i\}$  converges to a solution of the KS equation (5) with linear convergence rate no more than  $\max(1 - \alpha\gamma_{\max}^{-1}b_1, \alpha\gamma_{\min}^{-1}b_2 - 1)$ .

*Proof.* For any  $V^i$ , it follows from the definitions of  $D^i$ ,  $\alpha$  and  $b_2$  that

$$\begin{aligned} & \|I - \alpha(I - D^i)^{-1}(I - \partial_V \mathcal{V}(F_\phi(V^i)))\|_2 \\ &= \|I - \alpha(I - D^i) + \alpha(I - D^i)^{-1} \partial_V \mathcal{V}(F_\phi(V^i))\|_2 \\ &\leq \|I - \alpha(I - D^i)\|_2 + |\alpha| \|(I - D^i)^{-1} J(F_\phi(V^i)) J(V^i)\|_2 \\ &\leq \begin{cases} 1 - \alpha\gamma_{\max}^{-1} + \alpha\gamma_{\min}^{-1} \frac{\|L^\dagger\|_2 + \theta}{\delta} = 1 - \alpha\gamma_{\max}^{-1} b_1, & \text{if } \alpha < \gamma_{\max}; \\ \alpha\gamma_{\min}^{-1} - 1 + \alpha\gamma_{\min}^{-1} \frac{\|L^\dagger\|_2 + \theta}{\delta} = \alpha\gamma_{\min}^{-1} b_2 - 1, & \text{otherwise,} \end{cases} \\ &< 1. \end{aligned}$$

This completes the proof. □

## 5.1 Approximate Newton Method I

Our first approach replaces the operator  $\partial_V F_\phi(V^i)[\cdot]$  by a diagonal matrix  $\tau^i I$ , where  $\tau^i$  is a non-positive scalar. It is chosen to be non-positive since  $\partial_V F_\phi(V^i)[\cdot]$  is negative semidefinite from Lemma 4.10. Consequently, we set  $D^i := \tau^i J(\rho)$  and the scheme (106) becomes

$$V^{i+1} = V^i - \alpha (I - \tau^i J(F_\phi(V^i)))^{-1} (V^i - \mathcal{V}(F_\phi(V^i))). \quad (107)$$

The next theorem presents the local convergence analysis for the method (107).

**Theorem 5.2.** *Let  $V^*$  be a solution of the KS equation (5). Suppose that Assumption 4.1 holds with a constant  $\theta$  and Assumption 2.5 is valid at  $H(V^*)$  with a constant  $\delta$  satisfying*

$$\delta > -\lambda_{\min}^*, \quad (108)$$

where  $\lambda_{\min}^* := \min\{0, \lambda_{\min}(J(F_\phi(V^*)))\}$ . Let  $\{V^i\}$  be a sequence generated by the scheme (107) using  $\lim_{i \rightarrow \infty} \tau^i = \tau^* \in (-\frac{1}{\delta}, 0)$  and a step size

$$\alpha \in \left(0, \frac{\delta + \lambda_{\min}^*}{\|L^\dagger\|_2 + \theta + \delta}\right). \quad (109)$$

If the initial point  $V^0$  is selected in a sufficiently small open neighborhood of  $V^*$ , then  $\{V^i\}$  converges to  $V^*$  with  $R$ -linear convergence rate no more than

$$\max \left\{ \left( 1 - \alpha \left( \frac{\delta}{\|L^\dagger\|_2 + \theta + \delta} + \frac{\lambda_{\min}^*}{\delta + \lambda_{\min}^*} \right) \right), \left( \alpha \frac{\|L^\dagger\|_2 + \theta + \delta}{\delta + \lambda_{\min}^*} - 1 \right) \right\}.$$

*Proof.* The convergence of the iteration (107) is guaranteed by

$$\varrho(I - \alpha\mathcal{M}) < 1, \quad (110)$$

where  $\mathcal{M} = (I - \tau^* J(F_\phi(V^*)))^{-1} (I - J(F_\phi(V^*)) \partial_V F_\phi(V^*))$ . A direct linear algebraic calculation yields

$$\begin{aligned} \mathcal{M} &= (I - \tau^* J(F_\phi(V^*)))^{-1} - (I - \tau^* J(F_\phi(V^*)))^{-1} J(F_\phi(V^*)) \partial_V F_\phi(V^*) \\ &= I + (I - \tau^* J(F_\phi(V^*)))^{-1} J(F_\phi(V^*)) (\tau^* I - \partial_V F_\phi(V^*)). \end{aligned} \quad (111)$$

The symmetry of  $J(F_\phi(V^*))$  implies that  $(I - \tau^* J(F_\phi(V^*)))^{-1} J(F_\phi(V^*))$  is also symmetric, which together with the fact that  $\tau^* I - \partial_V F_\phi(V^*)$  is positive definite and Lemma 4.9 shows that all the eigenvalues of  $\mathcal{M}$  are real. Similar to the proof of Theorem 4.11, the inequality (110) holds if

$$\lambda_{\min}(\mathcal{M}) > 0; \quad (112)$$

$$\alpha \lambda_{\max}(\mathcal{M}) < 2. \quad (113)$$

Using  $0 > \tau^* > -\frac{1}{\delta}$  and the definition of  $\lambda_{\min}^*$ , we have

$$\lambda_{\min}(I - \tau^* J(F_\phi(V^*))) \geq \frac{\delta + \lambda_{\min}^*}{\delta} > 0, \quad (114)$$

$$\lambda_{\max}(I - \tau^* J(F_\phi(V^*))) \leq \frac{\|L^\dagger\|_2 + \theta + \delta}{\delta}. \quad (115)$$

Using the fact that the smallest eigenvalue of a summation of two matrices is larger than the summation of the smallest eigenvalues of these matrices, we obtain

$$\begin{aligned} \lambda_{\min}(\mathcal{M}) &\geq \lambda_{\min}((I - \tau^* J(F_\phi(V^*)))^{-1}) + \lambda_{\min}((I - \tau^* J(F_\phi(V^*)))^{-1} J(F_\phi(V^*)) (-\partial_V F_\phi(V^*))) \\ &\geq \frac{\delta}{\|L^\dagger\|_2 + \theta + \delta} + \lambda_{\min}((I - \tau^* J(F_\phi(V^*)))^{-1} J(F_\phi(V^*)) (-\partial_V F_\phi(V^*))). \end{aligned} \quad (116)$$

Applying Lemma 4.9,  $\lambda_{\max}(-\partial_V F_\phi(V^*)) \leq \frac{1}{\delta}$  from Lemma 4.2 and the definition of  $\lambda_{\min}^*$ , we have

$$\begin{aligned} &\lambda_{\min}((I - \tau^* J(F_\phi(V^*)))^{-1} J(F_\phi(V^*)) (-\partial_V F_\phi(V^*))) \\ &\geq \begin{cases} \lambda_{\min}((I - \tau^* J(F_\phi(V^*)))^{-1}) \lambda_{\min}(J(F_\phi(V^*))) \lambda_{\min}(-\partial_V F_\phi(V^*)), & \text{if } \lambda_{\min}(J(F_\phi(V^*))) \geq 0, \\ \lambda_{\max}((I - \tau^* J(F_\phi(V^*)))^{-1}) \lambda_{\min}(J(F_\phi(V^*))) \lambda_{\max}(-\partial_V F_\phi(V^*)), & \text{otherwise} \end{cases} \\ &\geq \begin{cases} 0, & \text{if } \lambda_{\min}(J(F_\phi(V^*))) \geq 0, \\ \frac{\lambda_{\min}(J(F_\phi(V^*)))}{\delta + \lambda_{\min}^*}, & \text{otherwise} \end{cases} \\ &\geq \frac{\lambda_{\min}^*}{\delta + \lambda_{\min}^*}, \end{aligned} \quad (117)$$

which together with (116) gives (112).

It follows from Lemma 4.9 and (114) that

$$\begin{aligned}\lambda_{\max}(\mathcal{M}) &\leq \lambda_{\max}((I - \tau^* J(F_\phi(V^*)))^{-1}) + \lambda_{\max}((I - \tau^* J(F_\phi(V^*)))^{-1} \lambda_{\max}(J(F_\phi(V^*))) \lambda_{\max}(-\partial_V F_\phi(V^*))) \\ &\leq \frac{\|L^\dagger\|_2 + \theta + \delta}{\delta + \lambda_{\min}^*}.\end{aligned}\quad (118)$$

Combining (109) and (118) together yields (113).  $\square$

Similar to Corollary (4.13), the condition (108) holds when  $J(F_\phi(V^*))$  is positive semidefinite.

## 5.2 Approximate Newton Method II

The matrix  $J(\rho)$  has to be calculated for each  $\rho$  in the approximate Newton method (107). If the computational cost of second-order derivatives of the exchange correlation function is expensive, a simpler choice is to approximate  $J(F_\phi(V^*))$  by  $L^\dagger$  and  $\partial_V F_\phi(V)$  by  $\tau^i I$ , that is,  $D^i = \tau^i L^\dagger$ . Hence, approximate Newton method (106) becomes

$$V^{i+1} = V^i - \alpha (I - \tau^i L^\dagger)^{-1} (V^i - \mathcal{V}(F_\phi(V^i))), \quad (119)$$

where  $\{\tau^i\}$  is negative. In fact, (119) is exactly the method of elliptic preconditioner proposed in [23].

**Theorem 5.3.** *Let  $V^*$  be a solution of the KS equation (5). Suppose that Assumption 4.1 holds with a constant  $\theta$  and Assumption 2.5 is valid at  $H(V^*)$  with a constant  $\delta$  satisfying*

$$\delta > \theta. \quad (120)$$

*Let  $\{V^i\}$  be a sequence generated by the scheme (107) using  $\lim_{i \rightarrow \infty} \tau_i = \tau^* \in \left(-\frac{1}{\xi}, 0\right)$  such that  $\xi \geq \frac{\|L^\dagger\|_2 \theta}{\delta - \theta}$ , and a step size*

$$\alpha \in \left(0, \frac{2}{\frac{\|L^\dagger\|_2 + \xi}{\xi} + \frac{\theta}{\delta}}\right). \quad (121)$$

*If the initial point  $V^0$  is selected in a sufficiently small open neighborhood of  $V^*$ , then  $\{V^i\}$  converges to  $V^*$  with  $R$ -linear convergence rate no more than*

$$\max \left\{ \left(1 - \alpha \left(\frac{\xi}{\|L^\dagger\|_2 + \xi} - \frac{\theta}{\delta}\right)\right), \left(\alpha \left(\frac{\|L^\dagger\|_2 + \xi}{\xi} + \frac{\theta}{\delta}\right) - 1\right) \right\}. \quad (122)$$

*Proof.* Let  $\bar{\mathcal{M}} = (I - \tau^* L^\dagger)^{-1} (I - \partial_V \mathcal{V}(F_\phi(V^*)))$ . The convergence of the iteration (119) is guaranteed by

$$\varrho(I - \alpha \bar{\mathcal{M}}) < 1. \quad (123)$$

Using the formulation of  $\partial_V \mathcal{V}(F_\phi(V^*))$ , we can decompose  $\bar{\mathcal{M}} = \bar{\mathcal{M}}_1 - \bar{\mathcal{M}}_2$ , where  $\bar{\mathcal{M}}_1 = (I - \tau^* L^\dagger)^{-1} (I - L^\dagger \partial_V F_\phi(V^*))$  and  $\bar{\mathcal{M}}_2 = (I - \tau^* L^\dagger)^{-1} (J(F_\phi(V^*)) - L^\dagger) \partial_V F_\phi(V^*)$ . Since  $L^\dagger$  is positive semidefinite, a similar proof as Theorem 5.2 implies that all the eigenvalues of  $\bar{\mathcal{M}}_1$  are real and

$$\lambda_{\min}(\mathcal{M}_1) > \frac{\xi}{\|L^\dagger\|_2 + \xi}, \quad (124)$$

$$\lambda_{\max}(\mathcal{M}_1) \leq \frac{\|L^\dagger\|_2 + \delta}{\delta}. \quad (125)$$

Using Assumption 4.1 and Lemma 4.2, we have

$$\begin{aligned} \|\bar{\mathcal{M}}_2\|_2 &= \|(I - \tau^* L^\dagger)^{-1} (J(F_\phi(V^*)) - L^\dagger) \partial_V F_\phi(V^*)\|_2 \\ &\leq \|(I - \tau^* L^\dagger)^{-1}\|_2 \|J(F_\phi(V^*)) - L^\dagger\|_2 \|\partial_V F_\phi(V^*)\|_2 \leq \frac{\theta}{\delta}. \end{aligned} \quad (126)$$

Using (124) and  $\xi \geq \frac{\|L^\dagger\|_2 \theta}{\delta - \theta}$ , we obtain

$$\lambda_{\min}(\bar{\mathcal{M}}_1) > \frac{\theta}{\delta}, \quad (127)$$

which together with (126) yields

$$(1 - \alpha \lambda_{\min}(\bar{\mathcal{M}}_1)) < 1 - \alpha \|\bar{\mathcal{M}}_2\|_2. \quad (128)$$

On the other hand, it follows from (121), (125) and (126) that

$$(\alpha \lambda_{\max}(\bar{\mathcal{M}}_1) - 1) < 1 - \alpha \|\bar{\mathcal{M}}_2\|_2. \quad (129)$$

Combining (128) and (129) together gives

$$\varrho(1 - \alpha \bar{\mathcal{M}}_1) < 1 - \alpha \|\bar{\mathcal{M}}_2\|_2. \quad (130)$$

which guarantees (123).  $\square$

## 6 Conclusion

The equivalence between the KS total energy minimization problem and the KS equation is ambiguous in the current literatures on KSDFT. A simple counter example shows that the solutions of these two problems are not necessarily the same. We examine the equivalence based on the optimality conditions for a specialized exchange correlation functional. We prove that a global solution of the KS minimization problem is a solution of the KS equation if the gap between the  $p$ th and  $(p + 1)$ st eigenvalues of the Hamiltonian  $H(X)$  is sufficiently large. The equivalence of a local minimizer requires that the corresponding charge densities are all positive. For strong local minimizers, the nonzero charge densities are bounded below by a positive constant uniformly. These properties are summarized in Table 1.

We improve the convergence analysis on the SCF iteration for solving the KS equation by analyzing the Jacobian of the corresponding fixed point maps. Global convergence of the simple mixing scheme can be established when there exists a gap between  $p$ th and  $(p + 1)$ st eigenvalues of the Hamiltonian  $H(X)$ . This assumption can be relaxed for local convergence analysis and if the charge density is computed using the Fermi-Dirac distribution. Our results requires much weaker conditions than the previous analysis in [24]. The structure of the Jacobian also suggests two approximate Newton methods. In particular, the second one is exactly the method of elliptic preconditioner proposed in [23]. Although our assumption on the gap is very stringent and is almost never satisfied in reality, our analysis is helpful for a better understanding of the KS minimization problem, the KS equation and the SCF iteration. A summary of our convergence results is presented in Table 2.

Table 1: Equivalence between the KS total energy minimization and the KS equation using the exchange correlation function  $e^T \epsilon_{xc}(\rho) = -\frac{3}{4}\gamma\rho^T\rho^{\frac{1}{3}}$

properties	eigenvalue gap $\delta$	Other Assumptions
A global minimizer $X^*$ solves the KS equation	Assumption 2.5 holds at $H(X^*)$ with $\delta > p (\ L^\dagger\ _2 - \frac{\gamma}{3})$	—
A local minimizer $X^*$ solves the KS equation	Assumption 2.5 holds at $H(X^*)$ with $\delta > 2 (\ L^\dagger\ _2 - \frac{\gamma}{3})$	$\rho_i > 0, i = 1, \dots, n$
$\rho_i(X^*) \in [0, c) \Rightarrow \rho_i(X^*) = 0$	—	$X^*$ is a strong local minimizer

## Acknowledgements

The authors would like to thank Dr. Chao Yang and Prof. Aihui Zhou for discussion on the KS equation and the SCF iteration.

## References

- [1] J. F. BONNANS AND A. SHAPIRO, *Perturbation Analysis of Optimization Problems*, Springer, New York, 2000.
- [2] C. LE BRIS, *Computational chemistry from the perspective of numerical analysis*, Acta Numer., 14 (2005), pp. 363–444.
- [3] ERIC CANCÈS, *SCF algorithms for Hartree-Fock electronic calculations*, Lecture Notes in Chemistry, 74 (2000), pp. 17–43.
- [4] ———, *Self-consistent field algorithms for Kohn-Sham models with fractional occupation numbers*, Journal of Chemical Physics, 114 (2001), p. 1061610622.
- [5] ERIC CANCÈS AND CLAUDE LE BRIS, *Can we outperform the DIIS approach for electronic structure calculations?*, International Journal of Quantum Chemistry, 79 (2000), pp. 82–90.
- [6] E. CANCÈS AND C. LE BRIS, *On the convergence of SCF algorithms for the Hartree-Fock equations*, Math. Model. Numer. Anal., 34 (2000), pp. 749–774.
- [7] E. CANCÈS, M. DEFRANCESCHI, W. KUTZELNIGG, C. LE BRIS, AND Y. MADAY, *Computational quantum chemistry: a primer*, in Handbook of numerical analysis. Volume X: special volume: computational chemistry, Ph. Ciarlet and C. Le Bris, eds., North-Holland, 2003, pp. 3–270.
- [8] ERIC CANCÈS AND KATARZYNA PERNAL, *Projected gradient algorithms for Hartree-Fock and density matrix functional theory calculations*, Journal of Chemical Physics, 128 (2008), pp. 108–134.
- [9] H. CHEN, X. DAI, X. GONG, L. HE, AND A. ZHOU, *Adaptive finite element approximations for Kohn-Sham models*. arXiv:1302.6896.
- [10] H. CHEN, X. GONG, L. HE, Z. YANG, AND A. ZHOU, *Numerical analysis of finite dimensional approximations of Kohn-Sham models*, Adv. Comput. Math., 38 (2013), pp. 225–256.

Table 2: convergence results for solving the KS equation under Assumption 4.1

properties		eigenvalue gap $\delta$ or smoothing parameter $\beta$	step size $\alpha$
analysis of SCF in [24]	global convergence	UWP holds and $\delta > 12p\sqrt{n}(\ L^\dagger\ _2 + \theta)$	1
	local convergence	Assumption 2.5 holds at the local minimizer with $\delta > 2\sqrt{n}(\ L^\dagger\ _2 + \theta)$	1
analysis of SCF with simple mixing	global convergence using $F_\phi$	UWP holds with $\delta > \ L^\dagger\ _2 + \theta$	$\left(0, \frac{2\delta}{\ L^\dagger\ _2 + \theta + \delta}\right)$
	global convergence using $F_{f_\mu}$	$\frac{4}{\beta} > \ L^\dagger\ _2 + \theta$	$\left(0, \frac{8}{(\ L^\dagger\ _2 + \theta)\beta + 4}\right)$
	local convergence using $F_\phi$	Assumption 2.5 holds at the local minimizer with $\delta > -\min\{0, \lambda_{\min}(J(F_\phi(V^*)))\}$	$\left(0, \frac{2\delta}{\ L^\dagger\ _2 + \theta + \delta}\right)$
	local convergence using $F_{f_\mu}$	$\frac{4}{\beta} > -\min\{0, \lambda_{\min}(J(F_\phi(V^*)))\}$	$\left(0, \frac{8}{(\ L^\dagger\ _2 + \theta)\beta + 4}\right)$
analysis of Approximate Newton methods	global convergence	UWP holds with $\delta > \frac{\gamma_{\max}}{\gamma_{\min}} \cdot (\ L^\dagger\ _2 + \theta)$	$\left(0, \frac{2\delta}{\ L^\dagger\ _2 + \theta + \delta}\right)$
	local convergence on $D^i := \tau^i J(\rho)$	Assumption 2.5 holds at the local minimizer with $\delta > -\min\{0, \lambda_{\min}(J(F_\phi(V^*)))\}$	$\left(0, \frac{\delta + \lambda_{\min}^*}{\ L^\dagger\ _2 + \theta + \delta}\right)$
	local convergence on $D^i := \tau^i L^\dagger$	Assumption 2.5 holds at the local minimizer with $\delta > \theta$	$\left(0, \frac{2\delta}{\frac{\delta}{\xi} \cdot (\ L^\dagger\ _2 + \xi) + \theta}\right)$

- [11] X. CHEN, H. QI, AND P. TSENG, *Analysis of nonsmooth symmetric-matrix-valued functions with applications to semidefinite complementarity problems*, SIAM Journal on Optimization, 13 (2003), pp. 960–985.
- [12] X. DAI, X. GONG, Z. YANG, D. ZHANG, AND A. ZHOU, *Finite volume discretizations for eigenvalue problems with applications to electronic structure calculations*, Multiscale Model. Simul., 9 (2011), pp. 208–240.
- [13] X. DAI, Z. YANG, AND A. ZHOU, *Symmetric finite volume schemes for eigenvalue problem in arbitrary dimensions*, Sci. China Ser. A., 51 (2008), pp. 1401–1414.
- [14] CHAO DING, *An Introduction to a Class of Matrix Optimization Problems*, PhD thesis, National University of Singapore, 2012.
- [15] D. DRUSVYATSKIY AND A. S. LEWIS, *Tilt stability, uniform quadratic growth, and strong metric regularity of the subdifferential*, SIAM Journal of Optimization, 23 (2013), pp. 256–267.
- [16] WEIGUO GAO, CHAO YANG, AND JUAN MEZA, *Solving a class of nonlinear eigenvalue problems by Newton’s method*, tech. report, Lawrence Berkeley National Laboratory, 2009.
- [17] G. P. KERKER, *Efficient iteration scheme for self-consistent pseudopotential calculations*, Phys. Rev. B, 23 (1981), pp. 3082–3084.
- [18] J. KOUTECKÝ AND V. BONACIC, *On the convergence difficulties in the iterative Hartree-Fock procedure*, J. Chem. Phys., 55 (1971), pp. 2408–2413.

- [19] G. KRESSE AND J. FURTHMULLER, *Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set*, Computational Materials Science, 6 (1996), pp. 15–50.
- [20] KONSTANTIN N. KUDIN, GUSTAVO E. SCUSERIA, AND ERIC CANCÈS, *A black-box self-consistent field convergence algorithm: One step closer*, Journal of Chemical Physics, 116 (2002), pp. 8255–8261.
- [21] ANTOINE LEVITT, *Convergence of gradient-based algorithms for the Hartree-Fock equations*, ESAIM: Mathematical Modelling and Numerical Analysis, 46 (2012), pp. 1321–1336.
- [22] A. LEWIS AND H. SENDOV, *Twice differentiable spectral functions*, SIAM Journal on Matrix Analysis and Applications, 23 (2001), pp. 368–386.
- [23] LIN LIN AND CHAO YANG, *Elliptic preconditioner for accelerating the self consistent field iteration in Kohn-Sham density functional theory*. SIAM Journal on Scientific Computing.
- [24] XIN LIU, XIAO WANG, ZAIWEN WEN, AND YAXIANG YUAN, *On the convergence of the self-consistent field iteration in Kohn-Sham density functional theory*, tech. report, 2013. arXiv:1302.6022.
- [25] JUDGE NOCEDAL AND STEPHEN WRIGHT, *Numerical Optimization*, Springer, 2006.
- [26] REINHOLD SCHNEIDER, THORSTEN ROHWEDDER, ALEXEY NEELOV, JOHANNES, AND BLAUERT, *Direct minimization for calculating invariant subspaces in density functional computations of the electronic structure*, Journal of Computational Mathematics, 27 (2009), pp. 360–393.
- [27] A. SHAPIRO, *On differentiability of symmetric matrix valued functions*. optimization online:2002/07/499.
- [28] M. TORKI, *Second-order directional derivatives of all eigenvalues of a symmetric matrix*, Nonlinear analysis, 46 (2001), pp. 1133–1150.
- [29] Z. WEN, A. MILZAREK, M. ULBRICH, AND H. ZHANG, *Adaptive regularized self-consistent field iteration with exact Hessian for electronic structure calculation*, SIAM Journal on Scientific Computing, 35 (2013), pp. A1299–A1324.
- [30] C. YANG, W. GAO, AND J. MEZA, *On the convergence of the self-consistent field iteration for a class of nonlinear eigenvalue problems*, SIAM J. Matrix Analysis Applications, 30 (2009), pp. 1773–1788.
- [31] LEI-HONG ZHANG AND REN-CANG LI, *Maximization of the sum of the trace ratio on the Stiefel manifold*, tech. report, Shanghai University of Finance and Economics, 2013.