

# Regularity and Well-posedness of a Dual Program for Convex Best $C^1$ -spline Interpolation <sup>\*</sup>

Houduo Qi <sup>†</sup> and Xiaoqi Yang <sup>‡</sup>

July 2, 2004; Revised August 18, 2005

**Abstract:** An efficient approach to computing the convex best  $C^1$ -spline interpolant to a given set of data is to solve an associated dual program by standard numerical methods (e.g., Newton's method.) We study regularity and well-posedness of the dual program: two important issues that have been not yet well-addressed in the literature. Our regularity results characterize the case when the generalized Hessian of the objective function is positive definite. We also give sufficient conditions for the coerciveness of the objective function. These results together specify conditions when the dual program is well-posed and hence justify why Newton's method is likely to be successful in practice. Examples are given to illustrate the obtained results.

**Key Words.** Convex best interpolation, splines, Newton method, regularity, well-posedness, degeneracy.

**AMS Subject Classification.** 41A29, 65D15, 49J52, 90C25

---

<sup>\*</sup>The work was supported by a grant from EPSRC for the first author and by a grant from the Research Grants Council of the Hong Kong Special Administrative Region, China (Project No. PolyU 5141/01E) for the second author.

<sup>†</sup>School of Mathematics, University of Southampton, Highfield, Southampton SO17 1BJ, UK. Email: [hdqi@soton.ac.uk](mailto:hdqi@soton.ac.uk).

<sup>‡</sup>Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong, China. Email: [mayangxq@polyu.edu.hk](mailto:mayangxq@polyu.edu.hk).

# 1 Introduction

Suppose we are given a set of data  $(x_i, y_i) \in \mathbb{R}^2$ ,  $i = 0, 1, \dots, n$ , satisfying the condition:

$$\begin{cases} \Delta : a = x_0 < x_1 < \dots < x_{n-1} < x_n = b \\ \tau_i := (y_i - y_{i-1})/(x_i - x_{i-1}), \quad \tau_{i+1} \geq \tau_i \quad i = 1, \dots, n-1. \end{cases} \quad (1)$$

Let  $V[a, b] := \text{Sp}(3, \Delta)$  denote the set of cubic  $C^1$  splines on  $\Delta$ . Specifically,  $V[a, b]$  consists of functions that are (once) continuously differentiable everywhere and the restrictions of those functions on each subinterval  $[x_{i-1}, x_i]$  are cubic.

The so-called convex best  $C^1$  spline interpolation over the given data set is the optimal solution  $s^*$  of the following minimization problem ([1]):

$$\begin{aligned} & \text{minimize} && \|s''\|_2 \\ & \text{subject to} && s(x_i) = y_i, \quad i = 0, 1, \dots, n \quad (\text{interpolation condition}) \\ & && s \text{ is convex on } [a, b] \quad (\text{shape constraint}) \\ & && s \in V[a, b] \quad (\text{function space}), \end{aligned} \quad (2)$$

where  $\|\cdot\|_2$  is the  $L_2$ -norm in the sense of Lebesgue<sup>1</sup>. Taking into account of the interpolation constraints  $s(x_i) = y_i$ , it is easy to see that  $s \in \text{Sp}(3, \Delta)$  if and only if the restriction of  $s$  on each subinterval  $[x_{i-1}, x_i]$  is in the following form:

$$\begin{aligned} s(x) = & y_{i-1} + m_{i-1}(x - x_{i-1}) \\ & + (3\tau_i - 2m_{i-1} - m_i) \frac{(x - x_{i-1})^2}{h_i} + (m_{i-1} + m_i - 2\tau_i) \frac{(x - x_{i-1})^3}{h_i^2}, \end{aligned} \quad (3)$$

with  $h_i := x_i - x_{i-1}$ ,  $i = 1, \dots, n$ , and  $m_{i-1}$  and  $m_i$  are parameters. Furthermore, the condition (1) is usually known to be a necessary condition for  $s$  being convex. A sufficient and necessary condition for  $s$  being convex is

$$2m_{i-1} + m_i \leq 3\tau_i \leq m_{i-1} + 2m_i, \quad i = 1, \dots, n, \quad (4)$$

see [14]. Therefore, the minimization problem (2) is equivalent to the following problem:

$$\begin{aligned} & \min && \|s''\|_2^2 \\ & \text{s.t.} && 2m_{i-1} + m_i \leq 3\tau_i \leq m_{i-1} + 2m_i, \quad i = 1, \dots, n, \end{aligned} \quad (5)$$

where  $s$  takes the form (3).

Derived either by the Karush-Kuhn-Tucker optimality theorem [1] or by the Fenchel conjugation theory [3, 21], the dual program of (5) is

$$\max_{p \in \mathbb{R}^{n-1}} -L(p) := -\sum_{i=1}^n \frac{h_i}{12} q(p_i, p_{i-1}) - \sum_{i=1}^{n-1} p_i (\tau_{i+1} - \tau_i), \quad (6)$$

---

<sup>1</sup>If the function space  $V[a, b]$  is chosen to be  $W^{2,2}[a, b]$  (the Sobolev space on  $[a, b]$ ), then (2) becomes the well known convex best interpolation problem. A historical account of this problem can be found in a recent book by Deutsch [2]. It has been recently revisited by Dontchev, Qi and Qi [5] in the similar spirit as the current paper. The  $W^{2,2}[a, b]$  case is much harder to solve than the  $\text{Sp}(3, \Delta)$  case, but they do not cover each other

where  $p_0 = p_n = 0$  and  $q : \Re^2 \rightarrow \Re$  is the continuously differentiable piecewise convex quadratic function given by

$$q(a, b) := \begin{cases} (a^2 + ab + b^2) & \text{for } a \leq 0, b \leq 0 \\ (\frac{1}{2}a + b)^2 & \text{for } a \geq 0, a + 2b \leq 0 \\ (a + \frac{1}{2}b)^2 & \text{for } b \geq 0, 2a + b \leq 0 \\ 0 & \text{for } a + 2b \geq 0, 2a + b \geq 0. \end{cases}$$

For easy reference, we write (6) in its equivalent minimization form:

$$\min_{p \in \Re^{n-1}} L(p) = \sum_{i=1}^n \frac{h_i}{12} q(p_i, p_{i-1}) + \sum_{i=1}^{n-1} p_i (\tau_{i+1} - \tau_i). \quad (7)$$

Since  $L$  is convex, the optimality condition of (7) has the form of nonlinear equations

$$F(p) = -d, \quad (8)$$

where  $d \in \Re^{n-1}$  with  $d_i = 12(\tau_{i+1} - \tau_i)$  and  $F : \Re^{n-1} \rightarrow \Re^{n-1}$  is given by

$$F_i(p) = h_{i+1} \partial_2 q(p_{i+1}, p_i) + h_i \partial_1 q(p_i, p_{i-1}), \quad i = 1, \dots, n-1. \quad (9)$$

Here,  $\partial_1 q$  and  $\partial_2 q$  are the partial derivatives of  $q$  with respect to its first and second arguments. Once a solution of (8), say  $\bar{p}$ , is obtained, the convex best  $C^1$  cubic spline which solves (2) can be constructed via (3) with  $m$  given by

$$\begin{aligned} m_{i-1} &= \tau_i - \frac{h_i}{12} \left( \bar{p}_i + \frac{1}{2} \bar{p}_{i-1}^+ - 2 \bar{p}_{i-1}^- \right)^-, \\ m_i &= \tau_i + \frac{h_i}{12} \left( \bar{p}_{i-1} + \frac{1}{2} \bar{p}_i^+ - 2 \bar{p}_i^- \right)^-, \end{aligned}$$

for  $i = 1, 2, \dots, n$ , where  $a^+ = \max(0, a)$  and  $a^- = -\min(0, a)$  for  $a \in \Re$  (see [1, (28)]).

Advantages of the dual approach include that the constrained problem (2) is solved by the unconstrained convex problem (7) and that the classical Newton matrix (Hessian of  $L$ ) is tridiagonal and positive semidefinite, but at the cost of that  $L$  is only once continuously differentiable (This implies that the Hessian of  $L$  may not exist for some points.) Despite this, numerical experiments [1, 18, 19, 20] show that the Newton method terminates in a small number of steps (averaging 3 – 5 steps) if a sufficiently good starting point is used. In [16], by making use of the piecewise linearity property of  $F$ , we explained that the lack of twice differentiability of  $L$  does not present any difficulty in analyzing the convergence of Newton's method. We achieved this by recasting the classical Newton method as a generalized Newton method, which generates the next iterate  $p^+$  from the current iterate  $p$  as follows

$$p^+ = p - V^{-1}(F(p) + d), \quad V \in \partial_B F(p), \quad (10)$$

where  $\partial_B F(p)$  is the  $B$ -differential of  $F$  at  $p$  [17]. When  $F$  is differentiable at  $p$ ,  $\partial_B F(p) = \{\nabla F(p)\}$  (see Lem. 2.1) and hence the generalized Newton method (10) becomes the classical Newton method. For this reason, (10) is also referred to as Newton's method. There

is a natural extension of the convergence result for the classical Newton method to the generalized Newton method. A result of Kojima and Shindo [12, Thm. 1] when applied to (10) says that it quadratically converges to its solution, say  $p^*$ , provided that the starting point is sufficiently close to  $p^*$  and every element in  $\partial_B F(p^*)$  is nonsingular. (In [17], this property of nonsingularity is called *BD-regularity*. Since  $L$  is convex, nonsingularity means that all the matrices in  $\partial_B F(p^*)$  are positive definite. This in turn implies that, the Clarke generalized Jacobian  $\partial F(p^*)$ , being the convex hull of  $\partial_B F(p^*)$ , also contains only positive definite matrices. For this reason, we say  $F$  regular if every element in  $\partial_B F(p^*)$  is nonsingular. For more discussion on various regularities, see [8]). In addition, Newton's method finds the solution in one step when the current iterate and the solution fall in the same piece where  $F$  is linear. This observation has also been used in [9, 11, 13, 22] in showing the finite termination of various Newton's methods for a number of problems, which can be reformulated as system of piecewise linear equations.

Two important issues concerning the success of numerical optimization methods [10, 15] for the problem (7) are its regularity and well-posedness. By regularity, we mean the regularity of  $F$ . The property is important, for example, for the Newton method (10) being well-defined and it is also sufficient for (10) or quasi-Newton methods [12] to be locally quadratically convergent. By well-posedness, we mean that (7) has a unique solution towards which every minimizing sequence converges. In [16], we established the following result on regularity: Let

$$W := \{(a, b) \mid a + 2b < 0 \text{ or } b + 2a < 0\}.$$

It was shown that  $F$  is regular at  $\bar{p} \in \mathbb{R}^{n-1}$  if

$$(\bar{p}_i, \bar{p}_{i-1}) \in W \quad \text{for } i = 1, \dots, n \text{ with } \bar{p}_0 = \bar{p}_n = 0. \quad (11)$$

In this paper we continue our efforts in establishing more theoretical results on regularity and well-posedness of (7) which support the use of Newton's method. Condition (11), termed as the nondegenerate condition in this paper, is extended to degenerate cases. Let  $\mathcal{I}(\bar{p})$  contain all degenerate indices  $i$  of  $\bar{p}$  such that  $(\bar{p}_i, \bar{p}_{i-1}) \notin W$  and  $|\mathcal{I}(\bar{p})|$  be the cardinality of  $\mathcal{I}(\bar{p})$ . We consider the regularity of  $F$  at a point  $\bar{p}$  when  $|\mathcal{I}(\bar{p})| \neq 0$  (note that condition (11) implies  $|\mathcal{I}(\bar{p})| = 0$ .) If  $|\mathcal{I}(\bar{p})| = 1$ , then  $F$  is regular at  $\bar{p}$ , see Prop. 3.6; and in the case  $|\mathcal{I}(\bar{p})| \geq 2$ ,  $F$  is regular at  $\bar{p}$  if between every two degenerate indices of  $\bar{p}$  there exists at least one strongly nondegenerate index, see Prop. 3.4. Both results are included in Sec. 3 and rely on more accurate description of matrix structure of all elements in  $\partial_B F(\bar{p})$ , see Sec. 2. In Sec. 4 we present a set of conditions which ensures the coerciveness of (7), see Prop. 4.1. An example is also constructed to show that violation of these conditions may lead to unboundedness of level sets of the function  $L$ .

## 2 Elementary results

In this section, we first present the definition of  $\partial_B$  for any local Lipschitz function and then characterize the matrix structure of  $\partial_B F$ . These results will facilitate our study on the regularity of  $F$  in Sec. 3

For any Lipschitz mapping  $G : \mathbb{R}^\ell \rightarrow \mathbb{R}^\ell$ , the  $B$ -differential of  $G$  at  $p \in \mathbb{R}^\ell$  is defined by

$$\partial_B G(p) := \left\{ \lim_{\substack{p^k \rightarrow p \\ p^k \in D_G}} \nabla G(p^k) \right\}, \quad (12)$$

where  $D_G$  denotes the set where  $G$  is differentiable and  $\nabla G(p)$  denotes its Jacobian at  $p \in D_G$ . The regularity of  $G$  at  $p$  means that every  $V \in \partial_B G(p)$  is nonsingular. It is known that  $\partial_B G(p)$  becomes singleton if  $G$  is strictly differentiable at  $p$ . Differentiability only is not enough for  $\partial_B G$  being singleton. For our special function  $F$  defined in (8), differentiability of  $F$  at a point means continuous differentiability of  $F$  at that point, see Lem. 2.1.

Now we turn our attention to the function  $F$  in (9). Since  $F(p) = \nabla L(p)$  and  $L$  is convex, every  $\nabla F(p)$ ,  $p \in D_F$  is symmetric and positive semidefinite. Therefore,  $\partial_B F(p)$  contains only symmetric and positive semidefinite matrices. Since each  $F_i$  depends only on  $(p_{i+1}, p_i, p_{i-1})$ , every element in  $\partial_B F(p)$  must be tri-diagonal. Because  $F$  is piecewise linear, each  $\nabla F(p)$ ,  $p \in D_F$  is a constant matrix. This means by (12) that, for any  $V \in \partial_B F(p)$  there exists  $p^k \in D_F$  sufficiently near to  $p$  such that  $V = \nabla F(p^k)$ . An important application of this observation is to study the regularity of  $F$  at any point  $p \in \mathbb{R}^{n-1}$ . In fact, we only need to consider the nonsingularity of  $\nabla F(p^k)$  for all nearby points,  $p^k \in D_F$ , of  $p$ .

To facilitate our analysis, we define several sets in  $\mathbb{R}^2$ :

$$\begin{aligned} W_+ &:= \{(a, b) \mid (a, b) \in W, ab \neq 0\} \\ W_- &:= \{(a, b) \mid (a, b) \in W, a < 0 \text{ and } b < 0\} \\ \bar{W} &:= \{(a, b) \mid a + 2b \geq 0 \text{ and } b + 2a \geq 0\} \\ \bar{W}_+ &:= \{(a, b) \mid a + 2b > 0 \text{ and } b + 2a > 0\} \\ R_0^{ab} &:= \{(a, b) \mid a + 2b = 0, a \geq 0\} \cup \{(a, b) \mid b + 2a = 0, b \geq 0\} \\ R_-^a &:= \{(a, b) \mid b = 0, a \leq 0\} \\ R_-^b &:= \{(a, b) \mid a = 0, b \leq 0\}. \end{aligned}$$

It is easy to see that  $W \cup \bar{W} = \mathbb{R}^2$ ,  $\bar{W}_+$  is the interior of  $\bar{W}$ ,  $R_0^{ab}$  is the boundary of  $\bar{W}$ , and  $W_+ = W \setminus (R_-^a \cup R_-^b)$ . We further define two functions  $f, g : \mathbb{R}^2 \rightarrow \mathbb{R}$  by

$$f(a, b) := \partial_2 q(a, b), \quad g(a, b) := \partial_1 q(a, b).$$

For any  $p \in \mathbb{R}^{n-1}$  with  $p_0 = p_n = 0$ , we recall that

$$F_i(p) = h_{i+1}f(p_{i+1}, p_i) + h_i g(p_i, p_{i-1}), \quad i = 1, \dots, n-1$$

and note that  $f(a, b)$  is not differentiable on  $R_0^{ab} \cup R_-^a$ ,  $g(a, b)$  is not differentiable on  $R_0^{ab} \cup R_-^b$ . We come to the following result concerning the differentiability of  $F$ .

**Lemma 2.1**  *$F$  is differentiable at  $p \in \mathbb{R}^{n-1}$  with  $p_0 = p_n = 0$  if and only if*

$$p_1 \neq 0, p_{n-1} \neq 0, \text{ and } (p_i, p_{i-1}) \in W_+ \cup \bar{W}_+, \quad i = 2, \dots, n-1. \quad (13)$$

*Moreover,  $F$  is continuously differentiable at  $p$  if and only if  $F$  is differentiable at  $p$ .*

**Proof.** We note that  $F$  is differentiable at  $p \in \mathfrak{R}^{n-1}$  if and only if all  $F_i$ ,  $i = 1, \dots, n-1$  are differentiable at  $p$ . This in turn implies that  $f$  and  $g$  are differentiable at all pairs  $(p_{i+1}, p_i)$ ,  $i = 1, \dots, n-2$  and  $f$  is differentiable at  $(0, p_{n-1})$  and  $g$  is differentiable at  $(p_1, 0)$  (note that  $p_0 = p_n = 0$ .) These conditions amount to (13). The analytical expression of  $L$  further implies that  $F$  is continuously differentiable at a point if and only if  $F$  is differentiable at that point.  $\square$

It then follows from Lem. 2.1 that  $D_F$  is given by

$$D_F = \{p \in \mathfrak{R}^{n-1} \mid p_1 \neq 0, p_{n-1} \neq 0, \text{ and } (p_i, p_{i-1}) \in W_+ \cup \bar{W}_+, i = 2, \dots, n-1\}$$

and  $\partial_B F(p) = \{\nabla F(p)\}$  for all  $p \in D_F$ . In the following we let  $\bar{p}$  denote a reference point and  $p$  denote points sufficiently near to  $\bar{p}$ .

**Lemma 2.2** *Let  $\bar{p} \in \mathfrak{R}^{n-1}$  and  $V \in \partial_B F(\bar{p})$ . Then  $V_{i,i-1}$ ,  $V_{i,i}$ , and  $V_{i,i+1}$  in the  $i$ th row ( $2 \leq i \leq n-2$ ) can only possibly take the following values*

$V_{i,i-1}$	$V_{i,i}$	$V_{i,i+1}$
0	0	0
$h_i$	$\frac{1}{2}h_i$	0
$h_i$	$2h_i$	0
0	$\frac{1}{2}h_{i+1}$	$h_{i+1}$
0	$2h_{i+1}$	$h_{i+1}$
$h_i$	$\frac{1}{2}(h_i + h_{i+1})$	$h_{i+1}$
$h_i$	$2(h_i + h_{i+1})$	$h_{i+1}$

and  $V_{1,1}, V_{1,2}$  in the first row and  $V_{n-1,n-2}, V_{n-1,n-1}$  in the last row take the following values

$V_{1,1}$	$V_{1,2}$	$V_{n-1,n-2}$	$V_{n-1,n-1}$
0	0	0	0
$\frac{1}{2}h_2$	$h_2$	$h_{n-1}$	$\frac{1}{2}h_{n-1}$
$2h_1$	0	0	$2h_n$
$2(h_1 + h_2)$	$h_2$	$h_{n-1}$	$2(h_{n-1} + h_n)$

**Proof.** We only need to prove the result for  $\bar{p} \in D_F$ , because for any  $V \in \partial_B F(\bar{p})$  there exists  $p \in D_F$  such that  $V = \nabla F(p)$ . Let  $2 \leq i \leq n-2$  be fixed and recall that

$$F_i(p) = h_{i+1}f(p_{i+1}, p_i) + h_i g(p_i, p_{i-1}).$$

According to (13),  $(\bar{p}_{i+1}, \bar{p}_i) \in W_+ \cup \bar{W}_+$  and  $(\bar{p}_i, \bar{p}_{i-1}) \in W_+ \cup \bar{W}_+$ . Hence there are four possible cases to be considered (note that both  $W_+$  and  $\bar{W}_+$  are open set in  $\mathfrak{R}^2$ ):

Case 1.  $(\bar{p}_{i+1}, \bar{p}_i) \in \bar{W}_+$  and  $(\bar{p}_i, \bar{p}_{i-1}) \in \bar{W}_+$ . For this case,  $q(p_{i+1}, p_i) = q(p_i, p_{i-1}) = 0$ , hence  $V_{i,i-1} = V_{i,i} = V_{i,i+1} = 0$ .

Case 2.  $(\bar{p}_{i+1}, \bar{p}_i) \in \bar{W}_+$  and  $(\bar{p}_i, \bar{p}_{i-1}) \in W_+$ . For this case,  $q(p_{i+1}, p_i) = 0$ , hence  $V_{i,i+1} = 0$ .  $(\bar{p}_i, \bar{p}_{i-1}) \in W_+$  could only have three subcases:

$$\begin{aligned} \text{If } p_i > 0, p_{i-1} < 0, & \text{ then } q(p_i, p_{i-1}) = (\tfrac{1}{2}p_i + p_{i-1})^2 & \implies V_{i,i-1} = h_i, V_{i,i} = \tfrac{1}{2}h_i. \\ \text{If } p_i < 0, p_{i-1} > 0, & \text{ then } q(p_i, p_{i-1}) = (\tfrac{1}{2}p_{i-1} + p_i)^2 & \implies V_{i,i-1} = h_i, V_{i,i} = 2h_i. \\ \text{If } p_i < 0, p_{i-1} < 0, & \text{ then } q(p_i, p_{i-1}) = (p_i + p_{i-1}p_i + p_{i-1}^2) & \implies V_{i,i-1} = h_i, V_{i,i} = 2h_i. \end{aligned}$$

Case 3.  $(\bar{p}_{i+1}, \bar{p}_i) \in W_+$  and  $(\bar{p}_i, \bar{p}_{i-1}) \in \bar{W}_+$ . For this case,  $q(p_i, p_{i-1}) = 0$ , hence  $V_{i,i-1} = 0$ .  $(\bar{p}_{i+1}, \bar{p}_i) \in W_+$  could only have three subcases:

$$\begin{aligned} \text{If } p_i > 0, p_{i+1} < 0, & \text{ then } q(p_{i+1}, p_i) = (\tfrac{1}{2}p_i + p_{i+1})^2 \implies V_{i,i} = \tfrac{1}{2}h_{i+1}, V_{i,i+1} = h_{i+1}. \\ \text{If } p_i < 0, p_{i+1} > 0, & \text{ then } q(p_{i+1}, p_i) = (\tfrac{1}{2}p_{i+1} + p_i)^2 \implies V_{i,i} = 2h_{i+1}, V_{i,i+1} = h_{i+1}. \\ \text{If } p_i < 0, p_{i+1} < 0, & \text{ then } q(p_i, p_{i-1}) = (p_i + p_{i+1}p_i + p_{i+1}^2) \implies V_{i,i} = 2h_{i+1}, V_{i,i+1} = h_{i+1}. \end{aligned}$$

Case 4.  $(\bar{p}_{i+1}, \bar{p}_i) \in W_+$  and  $(\bar{p}_i, \bar{p}_{i-1}) \in W_+$ . For this case, we could only have the following three subcases (they are combinations of subcases in Case 2 and Case 3).

$$\begin{aligned} \text{If } p_i > 0, p_{i-1} < 0 \text{ and } p_{i+1} < 0 & \implies V_{i,i-1} = h_i, V_{i,i} = \tfrac{1}{2}(h_i + h_{i+1}), V_{i,i+1} = h_{i+1}. \\ \text{If } p_i < 0, p_{i-1} > 0 \text{ and } p_{i+1} > 0 \text{ or } < 0 & \implies V_{i,i-1} = h_i, V_{i,i} = 2(h_i + h_{i+1}), V_{i,i+1} = h_{i+1}. \\ \text{If } p_i < 0, p_{i-1} < 0 \text{ and } p_{i+1} > 0 \text{ or } < 0 & \implies V_{i,i-1} = h_i, V_{i,i} = 2(h_i + h_{i+1}), V_{i,i+1} = h_{i+1}. \end{aligned}$$

This proves our results for  $2 \leq i \leq n-2$ . We note that for  $i = 1$  or  $n-1$ ,  $\bar{p}_0 = \bar{p}_n = 0$  by (13). We can prove this part similarly as above.  $\square$

Lem. 2.2 is helpful in calculating elements in  $\partial_B F(\bar{p})$ . Note that the  $i$ th row of  $V$  is uniquely determined by its diagonal element  $V_{i,i}$ . Let  $\bar{p} \in \mathfrak{R}^{n-1}$  be given and  $V \in \partial_B F(\bar{p})$ . Let  $\sigma_i$  be the coefficient of the  $i$ th diagonal element of  $V$ . Then according to Lem. 2.2,  $\sigma_i \in \{0, \frac{1}{2}, 2\}$ . The following result says more about the possible value of  $\sigma_i$ .

**Lemma 2.3** *Let  $\bar{p} \in \mathfrak{R}^{n-1}$  with  $\bar{p}_0 = \bar{p}_n = 0$  be given,  $V \in \partial_B F(\bar{p})$  and  $1 \leq i \leq n-1$ . If*

$$(\bar{p}_{i+1}, \bar{p}_i) \in W_+, \tag{14}$$

*then*

$$(\sigma_i, \sigma_{i+1}) \in \{(1/2, 2), (2, 1/2), (2, 2)\}.$$

**Proof.** Without loss of generality, we assume that  $\bar{p} \in D_F$  and  $V = \nabla F(\bar{p})$ . Recall that

$$\begin{aligned} F_i(p) &= h_{i+1}f(p_{i+1}, p_i) + h_i g(p_i, p_{i-1}), \text{ and} \\ F_{i+1}(p) &= h_{i+2}f(p_{i+2}, p_{i+1}) + h_{i+1}g(p_{i+1}, p_i). \end{aligned}$$

Case 1.  $\bar{p}_{i+1} < 0$ . Then using arguments in Lem. 2.2 we have

$$\begin{aligned} \text{If } \bar{p}_i > 0 & \text{ then } q(p_{i+1}, p_i) = (\tfrac{1}{2}p_i + p_{i+1})^2 \implies \sigma_i = 1/2, \sigma_{i+1} = 2. \\ \text{If } \bar{p}_i < 0 & \text{ then } q(p_{i+1}, p_i) = (p_i^2 + p_i p_{i+1} + p_{i+1}^2) \implies \sigma_i = \sigma_{i+1} = 2. \end{aligned}$$

Case 2.  $\bar{p}_{i+1} > 0$ . Since  $(\bar{p}_{i+1}, \bar{p}_i) \in W_+$ , we must have  $\bar{p}_i > 0$ . In this case,  $q(p_{i+1}, p_i) = (\tfrac{1}{2}p_i + p_{i+1})^2$ . Hence, using similar arguments in Lem. 2.2, we have  $\sigma_i = 1/2$  and  $\sigma_{i+1} = 2$ . This completes the proof.  $\square$

The result in Lem. 2.3 does not depend on whether or not  $(\bar{p}_i, \bar{p}_{i-1}) \in W_+$  or  $(\bar{p}_{i+2}, \bar{p}_{i+1}) \in W_+$ . However, it is not true if condition (14) is replaced by  $(\bar{p}_i, \bar{p}_{i-1}) \in W_+$ . Suppose this latter condition holds and  $(\bar{p}_{i+1}, \bar{p}_i) \in \bar{W}_+$ . In the case of  $\bar{p}_i > 0$ , we must have  $\bar{p}_{i-1} < 0$ . Hence  $\sigma_i = 1/2$ . If, furthermore,  $\bar{p}_{i+1} > 0$  and  $\bar{p}_{i+2} < 0$ , then  $\sigma_{i+1} = 1/2$ .

### 3 Regularity

To ease our discussion we first introduce some concepts of degenerate index, strongly degenerate index, nondegenerate index and strongly nondegenerate index of a given point.

**Definition 3.1** Let  $\bar{p} \in \mathfrak{R}^{n-1}$  be given with  $\bar{p}_0 = \bar{p}_n = 0$ . We call  $i \in \{1, \dots, n\}$

- (i) a degenerate index of  $\bar{p}$  if  $(\bar{p}_i, \bar{p}_{i-1}) \in \bar{W}$ ,
- (ii) a strongly degenerate index of  $\bar{p}$  if  $(\bar{p}_i, \bar{p}_{i-1}) \in \bar{W}_+$ ,
- (iii) a nondegenerate index of  $\bar{p}$  if  $(\bar{p}_i, \bar{p}_{i-1}) \in W$ , and
- (iv) a strongly nondegenerate index of  $\bar{p}$  if  $(\bar{p}_i, \bar{p}_{i-1}) \in W_-$ .

We further let  $\mathcal{I}(\bar{p})$  and  $\mathcal{I}_+(\bar{p})$  be the sets of all degenerate indices and strongly degenerate indices of  $\bar{p}$  respectively:

$$\begin{aligned}\mathcal{I}(\bar{p}) &:= \{i \mid (\bar{p}_i, \bar{p}_{i-1}) \in \bar{W}, i = 1, \dots, n\} = \{i_1, i_2, \dots, i_{r-1}\} \quad (1 \leq r-1 \leq n), \\ \mathcal{I}_+(\bar{p}) &:= \{i \mid (\bar{p}_i, \bar{p}_{i-1}) \in \bar{W}_+, i = 1, \dots, n\},\end{aligned}$$

and let  $i_0 = 1$ ,  $i_r = n$  and  $\kappa_j := i_j - i_{j-1}$ ,  $j = 1, \dots, r$ . Note that  $\mathcal{I}_+(\bar{p}) \subseteq \mathcal{I}(\bar{p})$ . It is possible that  $i_1 = i_0 = 1$  and/or  $i_r = i_{r-1} = n$ , implying  $\kappa_1 = \kappa_r = 0$ . With these numbers, we have the following result of characterizing the structure of elements in  $\partial_B F(\bar{p})$  when  $\mathcal{I}_+(\bar{p}) = \mathcal{I}(\bar{p})$  (all degenerate indices are strongly degenerate indices.)

**Lemma 3.2** Let  $\bar{p} \in \mathfrak{R}^{n-1}$  be given with  $\bar{p}_0 = \bar{p}_n = 0$  such that  $\mathcal{I}_+(\bar{p}) = \mathcal{I}(\bar{p})$ . Let the numbers  $\kappa_j$ ,  $j = 1, \dots, r$  be defined as above. Then  $F$  is continuously differentiable at  $\bar{p}$  and  $\partial_B F(\bar{p}) = \{\nabla F(\bar{p})\}$ . Moreover  $V := \nabla F(\bar{p})$  is a block diagonal matrix of size  $r$ , i.e.,

$$V = \begin{pmatrix} V_1 & & & \\ & V_2 & & \\ & & \ddots & \\ & & & V_r \end{pmatrix},$$

and each  $V_j$ ,  $j = 1, \dots, r$  is a  $\kappa_j \times \kappa_j$  tridiagonal matrix. Each block of  $V_j$  corresponds to the functions  $F_{i_{j-1}}, \dots, F_{i_j-1}$ .

**Proof.** The proof is easy. The condition  $\mathcal{I}_+(\bar{p}) = \mathcal{I}(\bar{p})$  means that (13) holds; equivalently  $F$  is continuously differentiable at  $\bar{p}$  and hence  $\partial_B F(\bar{p}) = \{\nabla F(\bar{p})\}$ . Recall that  $|\mathcal{I}(\bar{p})|$  is the cardinality of  $\mathcal{I}(\bar{p})$ . If  $|\mathcal{I}(\bar{p})| = 0$ , then there is only one diagonal block, which is  $V$  itself. Suppose now  $|\mathcal{I}(\bar{p})| \neq 0$  and let  $i_\ell \in \mathcal{I}(\bar{p})$ . We prove that  $V_{i_\ell, i_\ell-1} = 0$ , which implies  $V_{i_\ell-1, i_\ell} = 0$  by symmetry of  $V$ . Hence  $V$  is separated by block tridiagonal matrices at  $i_\ell$ th diagonal element. Note that

$$F_{i_\ell}(p) = h_{i_\ell+1}f(p_{i_\ell+1}, p_{i_\ell}) + h_{i_\ell}g(p_{i_\ell}, p_{i_\ell-1}).$$

Since  $(\bar{p}_{i_\ell}, \bar{p}_{i_\ell-1}) \in \bar{W}_+$  and  $\bar{W}_+$  is an open set,  $g(p_{i_\ell}, p_{i_\ell-1}) = 0$  for any  $p$  close to  $\bar{p}$ , that is  $V_{i_\ell, i_\ell-1} = 0$ . This proves the result.  $\square$

Lem. 3.2 means that  $V$  is positive definite if and only if each block  $V_j$  is positive definite. The following result focuses on all possible choices of  $2 \times 2$  blocks.



**Lemma 3.3** Let  $\bar{p} \in \mathbb{R}^{n-1}$  be given. Suppose there exists  $2 \leq \ell \leq n-3$  such that

$$(\bar{p}_\ell, \bar{p}_{\ell-1}) \in \bar{W}_+, \quad (\bar{p}_{\ell+1}, \bar{p}_\ell) \in W_+, \quad \text{and} \quad (\bar{p}_{\ell+2}, \bar{p}_{\ell+1}) \in \bar{W}_+.$$

Then, for any  $V \in \partial_B F(\bar{p})$ , we have

$$\begin{pmatrix} V_{\ell,\ell} & V_{\ell,\ell+1} \\ V_{\ell+1,\ell} & V_{\ell+1,\ell+1} \end{pmatrix} = \begin{pmatrix} 2h_{\ell+1} & h_{\ell+1} \\ h_{\ell+1} & \frac{1}{2}h_{\ell+1} \end{pmatrix} \quad \text{if } \bar{p}_{\ell+1} > 0$$

and

$$\begin{pmatrix} V_{\ell,\ell} & V_{\ell,\ell+1} \\ V_{\ell+1,\ell} & V_{\ell+1,\ell+1} \end{pmatrix} \in \left\{ \begin{pmatrix} \frac{1}{2}h_{\ell+1} & h_{\ell+1} \\ h_{\ell+1} & 2h_{\ell+1} \end{pmatrix}, \begin{pmatrix} 2h_{\ell+1} & h_{\ell+1} \\ h_{\ell+1} & 2h_{\ell+1} \end{pmatrix} \right\} \quad \text{if } \bar{p}_{\ell+1} < 0.$$

Hence  $V$  is not necessarily positive definite.

**Proof.** Since  $(\bar{p}_{\ell+1}, \bar{p}_\ell) \in W_+$ ,  $(\bar{p}_\ell, \bar{p}_{\ell-1}) \in \bar{W}_+$  and  $\bar{p}_{\ell+1} > 0$ , Case 3 in the proof of Lem. 2.2 with  $i = \ell$  implies  $V_{\ell,\ell} = 2h_{\ell+1}$ ,  $V_{\ell,\ell+1} = h_{\ell+1}$ . Since  $(\bar{p}_{\ell+2}, \bar{p}_{\ell+1}) \in \bar{W}_+$ ,  $(\bar{p}_{\ell+1}, \bar{p}_\ell) \in W_+$  and  $\bar{p}_{\ell+1} > 0$ , Case 2 in the proof of Lem. 2.2 with  $i = \ell + 1$  implies  $V_{\ell+1,\ell} = h_{\ell+1}$ ,  $V_{\ell+1,\ell+1} = \frac{1}{2}h_{\ell+1}$ . The case  $\bar{p}_{\ell+1} < 0$  can be proved similarly.  $\square$

A close look at the proof for the only positive definite matrix in Lem. 3.3 reveals that it occurs when  $(\bar{p}_{\ell+1}, \bar{p}_\ell) \in W_-$ , i.e.,  $\bar{p}_{\ell+1} < 0$  and  $\bar{p}_\ell < 0$ . This leads to the following main result in this section.

**Proposition 3.4** Let  $\bar{p} \in \mathbb{R}^{n-1}$  be given with  $\bar{p}_0 = \bar{p}_n = 0$ . Let

$$\mathcal{I}(\bar{p}) = \{i_1, \dots, i_{r-1}\} \quad \text{and} \quad i_0 = 1, \quad i_r = n.$$

Suppose that there exists  $\ell_j$  such that

$$i_{j-1} < \ell_j < i_j \quad \text{and} \quad (\bar{p}_{\ell_j}, \bar{p}_{\ell_j-1}) \in W_-, \quad \forall j = 1, \dots, r. \quad (15)$$

Then each  $V \in \partial_B F(\bar{p})$  positive definite.

**Proof.** Let  $V \in \partial_B F(\bar{p})$ , then there exists a sequence  $\{p^k\} \subset D_F$  converging to  $\bar{p}$  and  $\nabla F(p^k) \rightarrow V$ . Because of the piecewise linearity of  $F$  (i.e.,  $F$  is piecewise linear function with finitely many pieces),  $\{\nabla F(p^k)\}$  contains finitely many matrices. For simplicity, we drop the superscript of  $p^k$  and assume  $p$  is sufficiently close to  $\bar{p}$ . We will show that  $\nabla F(p)$  is positive definite under the condition of (15).

The case  $|\mathcal{I}(p)| = 0$  has already been dealt with in [16, Lem. 3.2]. Without loss of generality, we assume  $|\mathcal{I}(p)| \geq 1$ . Since  $\mathcal{I}(p) \subseteq \mathcal{I}(\bar{p})$  for all  $p$  sufficiently close to  $\bar{p}$ . Renumbering the indices in  $\mathcal{I}(p)$  if necessary (for simplicity, we are still ordering them as  $i_1, \dots, i_{r-1}$ ), we see that the condition (15) still hold for  $p$  because  $W_-$  is an open set and  $p$  is sufficiently close to  $\bar{p}$ .

According to Lem. 3.2,  $V$  consists of  $r$  tridiagonal blocks. Let  $V_j$  be one of them. It corresponds to functions  $F_{i_{j-1}}, \dots, F_{i_j-1}$ . Since  $i_{j-1}, i_j \in \mathcal{I}(p)$ , we must have

$$(p_{i_{j-1}}, p_{i_{j-1}-1}) \in \bar{W}_+, \quad (p_{i_{j-1}+1}, p_{i_{j-1}}) \in W_+, \dots, \quad (p_{i_j-1}, p_{i_j-2}) \in W_+, \quad \text{and} \quad (p_{i_j}, p_{i_j-1}) \in \bar{W}_+,$$

which means that all indices between every two consecutive degenerate indices  $i_{j-1}$  and  $i_j$  of  $p$  are nondegenerate. According to Lem. 2.3 and Lem. 2.2,  $V_j$  has the following structure

$$V_j = \begin{pmatrix} \sigma_{i_{j-1}} h_{i_{j-1}+1} & h_{i_{j-1}+1} & & \\ h_{i_{j-1}+1} & \ddots & \ddots & \\ & \ddots & \ddots & h_{i_j-1} \\ & & h_{i_j-1} & \sigma_{i_{j-1}} h_{i_j-1} \end{pmatrix}$$

with

$$V_{\ell,\ell} = \sigma_\ell(h_\ell + h_{\ell+1}), \quad \forall i_{j-1} + 1 \leq \ell \leq i_j - 2 \quad (16)$$

and

$$(\sigma_\ell, \sigma_{\ell+1}) \in \{(1/2, 2), (2, 1/2), (2, 2)\}, \quad \forall i_{j-1} \leq \ell \leq i_j - 2. \quad (17)$$

That is, it cannot occur that  $\sigma_\ell = \sigma_{\ell+1} = 1/2$  for some  $i_{j-1} \leq \ell \leq i_j - 2$ . Since there exists  $\ell_j$  such that (15) holds and noting that  $W_-$  is open and that

$$\begin{aligned} F_{\ell_j-1}(p) &= h_{\ell_j} f(p_{\ell_j}, p_{\ell_j-1}) + h_{\ell_j-1} g(p_{\ell_j-1}, p_{\ell_j-2}), \\ F_{\ell_j}(p) &= h_{\ell_j+1} f(p_{\ell_j+1}, p_{\ell_j}) + h_{\ell_j} g(p_{\ell_j}, p_{\ell_j-1}), \end{aligned}$$

we must have

$$\sigma_{\ell_j-1} = \sigma_{\ell_j} = 2. \quad (18)$$

Let  $u \in \Re^{n-1}$  and  $u = (u^1, \dots, u^r)^T$  with  $u^j = (u_{i_{j-1}}, \dots, u_{i_j-1})$ . It follows that

$$\begin{aligned} & u^j V_j (u^j)^T \\ &= \sum_{\ell=i_{j-1}}^{i_j-1} V_{\ell,\ell} u_\ell^2 + 2 \sum_{\ell=i_{j-1}}^{i_j-2} h_{\ell+1} u_\ell u_{\ell+1} \\ &= \sigma_{i_{j-1}} h_{i_{j-1}+1} u_{i_{j-1}}^2 + \sum_{\ell=i_{j-1}+1}^{i_j-2} \sigma_\ell (h_\ell + h_{\ell+1}) u_\ell^2 + \sigma_{i_j-1} h_{i_j-1} u_{i_j-1}^2 + 2 \sum_{\ell=i_{j-1}}^{i_j-2} h_{\ell+1} u_\ell u_{\ell+1} \\ &\geq h_{\ell_j} (u_{\ell_j-1}^2 + u_{\ell_j}^2) + \sum_{\ell=i_{j-1}}^{i_j-2} h_{\ell+1} \min \left( \frac{1}{2} u_\ell^2 + u_{\ell+1}^2 + 2u_\ell u_{\ell+1}, u_\ell^2 + \frac{1}{2} u_{\ell+1}^2 + 2u_\ell u_{\ell+1} \right) \\ &= h_{\ell_j} (u_{\ell_j-1}^2 + u_{\ell_j}^2) + \frac{1}{2} \sum_{\ell=i_{j-1}}^{i_j-2} h_{\ell+1} \min \left\{ (u_\ell + 2u_{\ell+1})^2, (2u_\ell + u_{\ell+1})^2 \right\}. \end{aligned} \quad (19)$$

The inequality uses the relations (16)–(18). Hence,  $u^j V_j (u^j)^T \geq 0$  and equality holds if and only if  $u_\ell = 0$  for all  $\ell = i_{j-1}, \dots, i_j - 1$ , i.e.,  $u^j = 0$ . Therefore,

$$u^T V u = \sum_{j=1}^r u^j V_j (u^j)^T \geq 0$$

and the equality holds if and only if  $u^j = 0$  for  $j = 1, \dots, r$ , i.e.,  $u = 0$ . This proves that  $\nabla F(p)$  is positive definite for all  $p$  sufficiently close to  $\bar{p}$  and  $p \in D_F$ . For any  $V \in \partial_B F(\bar{p})$  there must be a near point  $p \in D_F$  such that  $V = \nabla F(p)$ . Hence  $V$  is positive definite.  $\square$

We are now to make some comments on condition (15). Suppose  $d_i > 0$  for all  $i = 1, \dots, n-1$  (otherwise,  $d_i = 0$  means that  $\tau_{i+1} = \tau_i$  and that  $(x_{i+1}, y_{i+1}), (x_i, y_i)$  and  $(x_{i-1}, y_{i-1})$  are on one line and the best cubic spline on the interval  $[x_{i-1}, x_{i+1}]$  is the line itself), and let  $\bar{p}$  be a solution of (8). By the definition of  $F_i$ , it cannot happen that  $i, i+1 \in \mathcal{I}(\bar{p})$  for any  $i = 1, \dots, n$ . That is, between any two degenerate indices there exists at least one nondegenerate index of  $\bar{p}$ . If one of these nondegenerate indices happens to be strongly nondegenerate, then  $F$  is regular at  $\bar{p}$  according to Prop. 3.4. Hence, Newton's method can be used to find the solution. We also note that violation of condition (15) may lead to singular Jacobian of  $F$  at  $\bar{p}$ . The following example, taken from [16], illustrates such a possibility.

**Example 3.5** *The given data is:  $(x_0, y_0) = (0, 0)$ ,  $(x_1, y_1) = (1, 1)$ ,  $(x_2, y_2) = (2, 3)$ ,  $(x_3, y_3) = (3, 7)$ ,  $(x_4, y_4) = (4, 12)$  and  $(x_5, y_5) = (5, 27)$ . It is easy to verify that the unique solution is  $\bar{p} = (2, -13, 60, -42)$ .  $\mathcal{I}(\bar{p}) = \{1, 3\}$ . However, there are no strongly nondegenerate indices, i.e., condition (15) is not satisfied. It follows from Lem. 2.1 that  $F$  is continuously differentiable at the solution and hence  $\partial_B F(p^*)$  contains only one element (noticing all  $h_i = 1$ ):*

$$\begin{pmatrix} 1/2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ 0 & 0 & 1/2 & 1 \\ 0 & 0 & 1 & 4 \end{pmatrix}$$

*which is singular.*

Prop. 3.4 gives a sufficient condition for the regularity of  $F$  at a point where degeneracy is in presence. An extreme case of this result is when there is no degenerate indices, i.e.,  $|\mathcal{I}(\bar{p})| = 0$ . In this case, condition (15) means that there needs one strongly nondegenerate index to ensure the regularity. On the other hand, the result [16, Lem. 3.2] states that  $F$  is regular at  $\bar{p}$  if  $|\mathcal{I}(\bar{p})| = 0$ , no need of an extra strongly nondegenerate index to ensure the regularity. That is, in the case  $|\mathcal{I}(\bar{p})| = 0$ , Prop. 3.4 is stronger (not weaker) than [16, Lem. 3.2]. This motivate us to investigate if Prop. 3.4 can be strengthened when  $|\mathcal{I}(\bar{p})| \neq 0$ . In the next proposition, we will illustrate that this at least can be done when  $|\mathcal{I}(\bar{p})| = 1$ . The proposition can also be thought of a direct extension of [16, Lem. 3.2] from the nondegenerate case to the degenerate case.

**Proposition 3.6** *Let  $\bar{p} \in \mathbb{R}^{n-1}$  be given with  $\bar{p}_0 = \bar{p}_n = 0$ . Then  $F$  is regular at  $\bar{p}$  if  $|\mathcal{I}(\bar{p})| = 1$ .*

**Proof.** Denote that  $\mathcal{I}(\bar{p}) = \{\ell\}$ . As in the previous proofs, we only consider these points  $p$  which are sufficiently close to  $\bar{p}$  and  $p \in D_F$ . It is enough to prove that  $V := \nabla F(p)$  is positive definite for these  $p$ . Since  $\mathcal{I}(p) \subseteq \mathcal{I}(\bar{p})$  for all  $p$  sufficiently close to  $\bar{p}$ . There are two cases to be considered.  $\mathcal{I}(p) = \emptyset$  and  $\mathcal{I}(p) = \{\ell\}$ . In the former case, [16, Lem. 3.2] says that  $V$  is positive definite. In the latter case, we consider the following possibilities of  $\ell$ .

**Case 1.**  $\ell = 1$ . We must have  $p_1 > 0$  since  $(p_1, 0) \in \bar{W}_+$ . Then  $V$  is given by

$$V = \begin{pmatrix} \frac{1}{2}h_2 & h_2 & & \\ h_2 & \ddots & \ddots & \\ & \ddots & \ddots & h_{n-1} \\ & & h_{n-1} & 2(h_{n-1} + h_n) \end{pmatrix},$$

and  $V_{i,i} = \sigma_i(h_i + h_{i+1})$  for all  $i = 2, \dots, n-2$  and  $(\sigma_i, \sigma_{i+1}) \in \{(1/2, 2); (2, 1/2), (2, 2)\}$  and it cannot occur that  $\sigma_i = \sigma_{i+1} = 1/2$  for any  $2 \leq i \leq n-2$ . These properties on  $\sigma_i$  come from Lem. 2.3 because  $(p_{i+1}, p_i) \in W_+$  for all  $i = 1, \dots, n-1$ . Like (19) we must have for  $u \in \Re^{n-1}$  that

$$u^T V u \geq \frac{1}{2} \sum_{i=1}^{n-2} h_{i+1} \left\{ (u_i + 2u_{i+1})^2, (2u_i + u_{i+1})^2 \right\} + h_n u_{n-1}^2.$$

Hence,  $V$  is positive definite. The case that  $\ell = n-1$  can be proved similarly.

**Case 2.**  $1 < \ell < n-1$ . Then  $V$  is a block diagonal matrix of size 2, i.e.,

$$V = \begin{pmatrix} V_1 & \\ & V_2 \end{pmatrix}$$

and

$$V_1 = \begin{pmatrix} V_{1,1} & h_2 & & \\ h_2 & \ddots & \ddots & \\ & \ddots & \ddots & h_{\ell-1} \\ & & h_{\ell-1} & V_{\ell-1,\ell-1} \end{pmatrix} \text{ and } V_2 = \begin{pmatrix} V_{\ell,\ell} & h_{\ell+1} & & \\ h_{\ell+1} & \ddots & \ddots & \\ & \ddots & \ddots & h_{n-1} \\ & & h_{n-1} & V_{n-1,n-1} \end{pmatrix}$$

with

$$V_{1,1} = \begin{cases} 2(h_1 + h_2) & \text{if } \ell > 2 \\ 2h_1 & \text{if } \ell = 2, \end{cases} \quad V_{n-1,n-1} = \begin{cases} 2h_{n-1} & \text{if } \ell = n-2 \\ 2(h_{n-1} + h_n) & \text{if } \ell < n-2, \end{cases}$$

and  $V_{\ell-1,\ell-1} = \sigma_{\ell-1}h_{\ell-1}$ ,  $V_{\ell,\ell} = \sigma_{\ell}h_{\ell+1}$ ,  $V_{i,i} = \sigma_i(h_i + h_{i+1})$  for  $i \in \{2, \dots, n-2\}$  but  $i \neq \ell-1, \ell$ . These  $\sigma_i$  must satisfy  $(\sigma_i, \sigma_{i+1}) \in \{(1/2, 2), (2, 1/2), (2, 2)\}$  for  $i \neq \ell-1$ . It is easy to see that both  $V_1$  and  $V_2$  are positive definite. This finishes the proof that  $F$  is regular at  $\bar{p}$  if  $|\mathcal{I}(\bar{p})| = 1$ .  $\square$

The condition  $|\mathcal{I}(\bar{p})| = 1$  cannot be weakened to  $|\mathcal{I}(\bar{p})| \geq 2$ . A counter-example is Example 3.5 where  $|\mathcal{I}(\bar{p})| = 2$  and the Jacobian of  $F$  at  $\bar{p}$  is singular.

To summarize, we are able to make the following claims on the regularity of  $F$  at a point  $\bar{p}$ : (i)  $F$  is regular at  $\bar{p}$  if  $|\mathcal{I}(\bar{p})| \leq 1$  (Prop. 3.6 and [16, Lem. 3.2]); (ii) If  $|\mathcal{I}(\bar{p})| \geq 2$ ,  $F$  is regular at  $\bar{p}$  if between every two degenerate indices there exists at least one strongly nondegenerate index (Prop. 3.4).

## 4 Well-posedness

We recall that the function  $L$  is said to be coercive if  $L(p) \rightarrow +\infty$  as  $\|p\| \rightarrow +\infty$ . If  $L$  is coercive and  $F$  is regular at a solution  $p^*$ , then the minimization problem (7) is well-posed in the sense that it has a unique solution towards which every minimizing sequence converges. Well-posedness of problem (7) ensures that the steepest descent method can be used as the first stage to find a good enough starting point for Newton's method, which has been numerically observed efficient [1, 18, 18, 20]. In this section we present sufficient conditions for the coerciveness of  $L$ . The proof technique for the following result is motivated by [4, 6] where a more difficult problem than (2) is considered (i.e.,  $V[a, b] = W^{2,2}[a, b]$ , the Sobolev space.)

**Proposition 4.1** *Suppose that  $d_i > 0$  for all  $i = 1, \dots, n-1$  and one of the following two conditions holds:*

$$2d_{i-1} > d_i, \quad \forall i = 2, \dots, n-1, \quad (20)$$

$$2d_{i+1} > d_i, \quad \forall i = 1, \dots, n-2. \quad (21)$$

*Then  $L$  is coercive.*

**Proof.** It is sufficient to prove that the level set

$$\mathcal{L}(c) := \{p \in \mathbb{R}^{n-1} \mid L(p) \leq c\}$$

is bounded for every  $c \in \mathbb{R}$ . Note that, for every  $c \in \mathbb{R}$ , the set  $\mathcal{L}(c)$  is closed and convex. Assume, on the contrary, that  $\mathcal{L}(c_0)$  is unbounded for some  $c_0 \in \mathbb{R}$  and let, without loss of generality,  $c_0 > 0$ . We first show that there exists a vector  $p \in \mathbb{R}^{n-1}$ ,  $p \neq 0$  such that  $\beta p \in \mathcal{L}(c_0)$  for every  $\beta \geq 0$ . Suppose that for every  $p \in \mathbb{R}^{n-1}$  there exists  $\beta_p \geq 0$  such that  $\beta_p p \notin \mathcal{L}(c_0)$ . From the convexity of  $\mathcal{L}(c_0)$  and  $0 \in \mathcal{L}(c_0)$ , it follows that  $\beta p \notin \mathcal{L}(c_0)$  whenever  $\beta \geq \beta_p$ . Let

$$\beta(p) := \max\{\beta \mid \beta \geq 0, \beta p \in \mathcal{L}(c_0)\}.$$

Then  $\beta(p) < +\infty$  since  $\mathcal{L}(c_0)$  is closed. It is also known from Dontchev and Kalchev [4, P.675] that  $\beta(\cdot)$  is an upper semicontinuous function over  $\mathbb{R}^{n-1}$ . Then

$$\beta^* := \sup\{\beta(p) \mid \|p\| = 1\} < +\infty.$$

Hence  $\mathcal{L}(c_0)$  is contained in a ball centered at the origin with radius  $\beta^* + 1$ . This contradiction establishes the existence of a vector  $p \in \mathbb{R}^{n-1}$ ,  $p \neq 0$  such that  $\beta p \in \mathcal{L}(c_0)$  for all  $\beta \geq 0$ . Now, for such  $p$ , we define

$$\kappa(\beta) := L(\beta p) = \frac{\beta^2}{12} \sum_{i=1}^n h_i q(p_i, p_{i-1}) + \frac{\beta}{12} \sum_{i=1}^{n-1} p_i d_i,$$

where we let  $p_0 = p_n = 0$ .

First we note that  $q(p_i, p_{i-1}) \geq 0$  for all  $i = 1, \dots, n$ . If there exists  $i_0 \in \{1, \dots, n\}$  such that  $(p_{i_0}, p_{i_0-1}) \in W$  then  $q(p_{i_0}, p_{i_0-1}) > 0$  and

$$\kappa(\beta) \geq \frac{\beta}{12} \left( \beta h_{i_0} q(p_{i_0}, p_{i_0-1}) - \sum_{i=1}^{n-1} p_i d_i \right) \rightarrow +\infty \quad \text{as } \beta \rightarrow +\infty.$$

This contradicts  $\beta p \in \mathcal{L}(c_0)$  for all  $\beta \geq 0$  (i.e.,  $\kappa(\beta) \leq c_0$ ). Hence, no pairs  $(p_i, p_{i-1}) \in W$  for any  $i = 1, \dots, n$ , resulting in

$$\kappa(\beta) = \frac{\beta}{12} \sum_{i=1}^{n-1} p_i d_i.$$

In the following we prove  $\sum_{i=1}^{n-1} p_i d_i > 0$ . We note once again that  $p \neq 0$ ,  $p_0 = p_n = 0$  and  $(p_i, p_{i-1}) \notin W$  for all  $i = 1, \dots, n$ . Define

$$\mathcal{J} := \{j \mid p_j < 0, \quad j = 1, \dots, n-1\}.$$

If  $\mathcal{J} = \emptyset$  it is obvious that  $\sum_{i=1}^{n-1} p_i d_i > 0$ . If  $j \in \mathcal{J}$  then  $p_{j-1} > 0$  and  $p_{j+1} > 0$  (i.e.,  $j-1, j+1 \notin \mathcal{J}$ ) because  $(p_i, p_{i-1}) \notin W$  for  $i = j, j+1$ . This yields

$$p_{j-1} + 2p_j \geq 0 \quad \text{and} \quad p_{j+1} + 2p_j \geq 0, \quad \forall j \in \mathcal{J}.$$

Hence

$$\begin{aligned} \sum_{i=1}^{n-1} p_i d_i &\geq \sum_{j \in \mathcal{J}} (p_j d_j + p_{j-1} d_{j-1}) \\ &\geq \sum_{j \in \mathcal{J}} (p_j d_j - 2p_j d_{j-1}) \\ &= \sum_{j \in \mathcal{J}} (2d_{j-1} - d_j)(-p_j) > 0 \quad (\text{because of (20)}). \end{aligned} \tag{22}$$

We call (22) backward grouping of  $(p_j, p_{j-1})$  over  $j \in \mathcal{J}$ . Similarly, we have forward grouping

$$\begin{aligned} \sum_{i=1}^{n-1} p_i d_i &\geq \sum_{j \in \mathcal{J}} (p_j d_j + p_{j+1} d_{j+1}) \\ &\geq \sum_{j \in \mathcal{J}} (p_j d_j - 2p_j d_{j+1}) \\ &= \sum_{j \in \mathcal{J}} (2d_{j+1} - d_j)(-p_j) > 0 \quad (\text{because of (21)}). \end{aligned}$$

This proves  $\sum_{i=1}^{n-1} p_i d_i > 0$ . Hence,  $\kappa(\beta) \rightarrow +\infty$  as  $\beta \rightarrow +\infty$ , a contradiction to  $\kappa(\beta) \leq c_0$ . Therefore,  $\mathcal{L}(c)$  is bounded for every  $c \in \mathfrak{R}$ . This finishes the proof.  $\square$

The following example shows that violation of conditions in Prop. 4.1 may lead to unboundedness of the level set  $\mathcal{L}(c_0)$  for some  $c_0 \in \mathfrak{R}$ .

**Example 4.2** Suppose we have data:  $(x_0, y_0) = (0, 0)$ ,  $(x_1, y_1) = (1, 1)$ ,  $(x_2, y_2) = (2, 3)$ ,  $(x_3, y_3) = (3, 9)$ ,  $(x_4, y_4) = (4, 16)$ . A solution to this problem is:  $p^* = (0, -12, 0)$ . Then

$$(d_1, d_2, d_3) = (12, 48, 12).$$

Let  $\bar{p} \in \mathfrak{R}^3$  be

$$(\bar{p}_1, \bar{p}_2, \bar{p}_3) = (2, -1, 2).$$

It is easy to see that  $|\mathcal{I}(\bar{p})| = 4$ . Hence

$$L(\beta\bar{p}) = \beta(\bar{p}_1d_1 + \bar{p}_2d_2 + \bar{p}_3d_3) = 0.$$

It follows that

$$\beta\bar{p} \in \mathcal{L}(0), \quad \forall \beta \geq 0$$

and  $\mathcal{L}(0)$  is unbounded. We see that condition (20) is not satisfied at  $(d_1, d_2)$  and condition (21) is not satisfied at  $(d_2, d_3)$ .

We conclude this section by two remarks

- (i) The minimization problem (7) is well-posed if  $L$  is coercive and  $F$  is regular at a solution. This paper studies when the coerciveness and regularity hold.
- (ii) Our regularity results have an algorithmic implication. Let  $p$  be the current iterate. If  $|\mathcal{I}(p)| \leq 1$  or  $|\mathcal{I}(p)| \geq 2$  and condition (15) is satisfied (with  $\bar{p}$  replaced by  $p$ ), then every  $V \in \partial_B F(p)$  is positive definite. Newton's method (10) successfully generates the next iterate  $p^+$ . If  $|\mathcal{I}(p)| \geq 2$  but condition (15) is not satisfied, then  $V \in \partial_B F(p)$  is possibly singular. In this case, we suggest to use the damped Newton method:

$$p^+ := p - (V + \epsilon I)^{-1}(F(p) + d),$$

where  $\epsilon > 0$  is a small positive number and  $I$  is the identity matrix. An alternative is to use the smoothing Newton method, as pointed out to us by one of the referees. The linear system solved in each iteration of the method is well defined. It remains to see how it performs compared to the Newton's method.

**Acknowledgement.** We thank the three referees for their detailed comments which greatly improved the presentation of the paper. Especially, a referee clarifies the use of the well-posedness.

## References

- [1] W. BURMEISTER, W. HESS, AND J.W. SCHMIDT, *Convex spline interpolants with minimal curvature*, Computing 35 (1985), 219–229.
- [2] F. DEUTSCH, *Best Approximation in Inner Product Spaces*, CMS Books in Mathematics 7. Springer-Verlag, New York, 2001.
- [3] S. DIETZE AND J.W. SCHMIDT, *Determination of shape preserving spline interpolants with minimal curvature via dual programs*, J. Approx. Theory 52 (1988), pp. 43–57.
- [4] A. L. DONTCHEV AND B. D. KALCHEV, *Duality and well-posedness in convex interpolation*, Numer. Funct. Anal. and Optim. 10 (1989), pp. 673–689.
- [5] A. L. DONTCHEV, H.-D. QI, AND L. QI, *Convergence of Newton's method for convex best interpolation*, Numer. Math. 87 (2001), pp. 435–456.

- [6] A. L. DONTCHEV, H.-D. QI, L. QI, AND H. YIN, *A Newton method for shape-preserving spline interpolation*, SIAM J. Optim. 13 (2002), pp. 588–602.
- [7] A.L. DONTCHEV AND T. ZOLEZZI, *Well-posed Optimization Problems*, Lecture Notes in Mathematics, 1543. Springer-Verlag, Berlin, 1993.
- [8] F. FACCHINEI, A. FISCHER, AND C. KANZOW, *Regularity properties of a new equation reformulation of variational inequalities*, SIAM J. Optim. 8 (1998), pp. 850–869.
- [9] A. FISCHER AND C. KANZOW, *On finite termination of an iterative method for linear complementarity problems*, Math. Programming 74 (1996), pp. 279–292.
- [10] J-B HIRIART-URRUTY AND C. LEMARÉCHAL *Convex Analysis and Minimization Algorithms I*, Springer-Verlag 1993.
- [11] C. KANZOW, H.-D. QI, AND L. QI, *On the minimum norm solution of linear programs*, J. Optim. Theory and Appl. 116 (2003), pp. 333–345.
- [12] M. KOJIMA AND S. SHINDO, *Extension of Newton and quasi-Newton methods for systems of  $PC^1$  equations*, J. Oper. Res. Soc. Japan 29 (1986), pp. 352–373.
- [13] O.L. MANGASARIAN, *Finite Newton method for classification problems*, Optim. Method and Software 17 (2002), pp. 913–929.
- [14] E. NEUMAN, *Uniform approximation by some Hermite interpolation splines*, J. Comput. Appl. Math. 4 (1978), pp. 7–9.
- [15] J. NOCEDAL AND S.J. WRIGHT, *Numerical Optimization*, Springer-Verlag, 1999.
- [16] H.-D. QI AND L. QI, *Finite termination of a dual Newton method for convex best  $C^1$  interpolation and smoothing*, Numer. Math., 96 (2003), pp. 317–337.
- [17] L. QI, *Convergence analysis of some algorithms for solving nonsmooth equations*, Math. Oper. Res. 18 (1993), pp. 227–244.
- [18] J.W. SCHMIDT, *An unconstrained dual program for computing convex  $C^1$ -spline approximants*, Computing 39 (1987), pp. 133–140.
- [19] J.W. SCHMIDT, *On tridiagonal linear complementarity problems*, Numer. Math. 51 (1987), pp. 11–21.
- [20] J.W. SCHMIDT AND W. HESS, *Spline interpolation under two-sided restrictions on the derivatives*, Z. Angew. Math. Mech. 69 (1989), pp. 353–365.
- [21] J.W. SCHMIDT, *Dual algorithms for solving convex partially separable optimization problems*, Jber. d. Dt. Math.-Verein. 94 (1992), 40–62.
- [22] D. SUN, J. HAN, AND Y. ZHAO, *The finite termination of the damped Newton algorithm for linear complementarity problems*, (Chinese) Acta Math. Appl. Sinica 21 (1998), pp. 148–154.