

decrease more in south town. And the plot $\alpha_{13}(\text{AGE})$ also indicates that a newer house is more expensive than a older one with similar condition.

Table 3.4 reports the averaged R^2 , MAEE, MAPE and the percentage of ‘good’ predicted price over 10 replicates for SCAD, FULL and parametric model respectively. From Table 3.4, FULL gives a slightly better estimation than SCAD at the expense of a more complex model structure. However for prediction, SCAD outperforms FULL. Therefore, SCAD not only gives a simpler and more interpretable model, but also improves prediction accuracy. Furthermore, the parametric model gives the worst results in both estimation and prediction, indicating that the data contains a nonlinear structure that can not be fully explained by the parametric model. Figure 3.4 plots the randomly selected actual prices against corresponding predicted prices under three models separately, with criteria band enclosing ‘good’ estimates. Again, it suggests both penalize model and full model do much better than the parametric model.

3.7 Proof of Lemmas and Theorems

Based on the oracle estimator $\hat{\gamma}^{(0)} = (\hat{\gamma}_{ls}^{(0)}, l = 1, \dots, d_1, s = 0, \dots, d_2)$ in Section 3.3, we further define $\hat{\alpha}_{ls}^{(0)}$ for notation convenience. For $(l, 0) \in S^{(0)}$, $\hat{\alpha}_{l0}^{(0)} = \hat{\gamma}_{l0}^{(0)}$. For $s > 0$ and $(l, s) \in S^{(0)}$, $\hat{\alpha}_{ls}^{(0)} = \hat{\gamma}_{ls}^{(0)} \mathbf{B}_s$. Here \mathbf{B}_s is the vector of the empirically centered spline basis on x_s defined in Section 3.2. For $(l, s) \notin S^{(0)}$, $\hat{\alpha}_{ls}^{(0)} = 0$.

For any square matrix \mathbf{U} , we denote $\rho_{min}(\mathbf{U})$ and $\rho_{max}(\mathbf{U})$ as the minimal and maximal eigenvalues of \mathbf{U} respectively. For notation simplicity, we use the same c, c_1, c_2 as general notations for positive constants with not necessarily the same value.

3.7.1 Preliminary Lemmas

Lemma 3.7.1. *Under Assumptions (C.3) - (C.4), for each pair of $(l, s) \in S$, the eigenvalues of \mathbf{W}_{ls} are bounded by two positive constants with probability approaching to 1.*

That is, let $\rho_{\min}(\mathbf{W}_{ls})$ and $\rho_{\max}(\mathbf{W}_{ls})$ be the minimal and maximal eigenvalues of \mathbf{W}_{ls} respectively. Then there exist $0 < c_1 < c_2$, such that

$$P(c_1 \leq \rho_{\min}(\mathbf{W}_{ls}) \leq \rho_{\max}(\mathbf{W}_{ls}) \leq c_2) \rightarrow 1 \text{ as } n \rightarrow \infty. \quad (3.8)$$

Proof. From Lemma 5 of Xue and Qu (2012), we know that (3.8) is true for $(l, s) \in S$ and $s = 0$. Now, for $s = 0$, $\rho(\mathbf{W}_{l0}) = \rho\left(\frac{\mathbf{t}_l^T \mathbf{t}_l}{n}\right) = \|T_l\|_n^2$. Therefore, according to Assumption (C.4), one has that $\|T_l\|_n^2 \leq c^2$. Furthermore, the Central Limit Theorem gives that, $\|T_l\|_n^2 \geq E(T_l^2) - o_P(1) \geq \text{Var}(T_l)$. Therefore, by letting $c_1 = \min_{l=1, \dots, d_1} \text{Var}(T_l)$ and $c_2 = c^2$, Lemma 3.7.1 is proved.

Lemma 3.7.2. Let $\mathbf{Z} = (\mathbf{Z}_{ls}, (l, s) \in S)$, $\mathbf{W} = \frac{\mathbf{Z}^T \mathbf{Z}}{n}$, then under Assumptions (C.3) - (C.4), the eigenvalues of \mathbf{W} are bounded within two positive constants. That is, there exist $c_2 > c_1 > 0$ such that, except on an event whose probability tends to zero, as $n \rightarrow \infty$,

$$c_1 \leq \rho_{\min}(\mathbf{W}) \leq \rho_{\max}(\mathbf{W}) \leq c_2. \quad (3.9)$$

Proof. For any coefficients $\gamma = (\gamma_{ls}^T, (l, s) \in S)^T \in \mathbb{R}^{d_1(1+d_2J_n)}$, one has

$$\gamma^T \mathbf{W} \gamma = \sum_{l=1}^{d_1} \left(\gamma_{l0} + \sum_{s=1}^{d_2} \sum_{j=1}^{J_n} \gamma_{ls,j} B_{s,j} T_l \right)_n^2.$$

Lemma A.5 in Xue and Yang (2006b) ensures that there exists a positive constants c_1 such that $\gamma^T \mathbf{W} \gamma \geq c_1 \sum_{l=1}^{d_1} \gamma_{l0}^2 + \sum_{s=1}^{d_2} \sum_{j=1}^{J_n} \gamma_{ls,j}^2 = c_1 \|\gamma\|^2$. On the other hand, Cauchy-Schwarz inequality gives that there exists a constant $c > 0$ such that

$$\gamma^T \mathbf{W} \gamma = \sum_{l=1}^{d_1} \left(\gamma_{l0} T_l + \sum_{s=1}^{d_2} \sum_{j=1}^{J_n} \gamma_{ls,j} B_{s,j} T_l \right)_n^2 \leq c \sum_{l=1}^{d_1} \gamma_{l0}^T \frac{\mathbf{Z}_{l0}^T \mathbf{Z}_{l0}}{n} \gamma_{l0} + \sum_{s=1}^{d_2} \gamma_{ls}^T \frac{\mathbf{Z}_{ls}^T \mathbf{Z}_{ls}}{n} \gamma_{ls}.$$

Therefore, from Lemma 3.7.1, there exists a positive c_2 such that $\gamma^T \mathbf{W} \gamma \leq c_2 \|\gamma\|^2$.

As a consequence of Lemmas 3.7.1 and 3.7.2, one has the following corollary.

Corollary 3.7.1.1. Let A be a subset of the full index pair set S . Denote $\mathbf{Z}_A = (\mathbf{Z}_{ls}, (l, s) \in A)$ and $\mathbf{W}_A = \frac{\mathbf{Z}_A^T \mathbf{Z}_A}{n}$, then under Assumption (C.3) - (C.4), there exist two positive

constants c_1, c_2 , such that

$$P(c_1 \leq \rho_{\min}(\mathbf{W}_A) \leq \rho_{\max}(\mathbf{W}_A) \leq c_2) \rightarrow 1 \text{ as } n \rightarrow +\infty.$$

Lemma 3.7.3. *Under Assumptions (C.3) - (C.4) and (C.6), each additive term of oracle estimators, $\hat{\alpha}_{l_0}^{(0)}$, converges to the corresponding true function α_{l_0} in probability. Specifically, we have*

$$\sum_{l=1}^{d_1} \hat{\alpha}_{l_0}^{(0)} T_l - \alpha_{l_0} T_l + \sum_{l=1}^{d_1} \sum_{s=1}^{d_2} \hat{\alpha}_{l_s}^{(0)}(X_s) T_l - \alpha_{l_s}(X_s) T_l = O_P \left(N_n^{-(p+1)} + \sqrt{\frac{N_n}{n}} \right).$$

Proof. Theorem 1 in Xue and Yang (2006) entails that

$$\max_{1 \leq l \leq d_1} \hat{\alpha}_{l_0}^{(0)} - \alpha_{l_0} + \max_{1 \leq l \leq d_1, 1 \leq s \leq d_2} \hat{\alpha}_{l_s}^{(0)}(X_s) - \alpha_{l_s}(X_s) = O_P \left(N_n^{-(p+1)} + \sqrt{\frac{N_n}{n}} \right).$$

Therefore, Lemma follows from Assumption (C.4).

3.7.2 Proof of Theorem 3.3.1

For notation simplicity, denote

$$L_n(\gamma) = \frac{1}{2n} \mathbf{Y} - \sum_{l=1}^{d_1} \sum_{s=0}^{d_2} \mathbf{Z}_{l_s}(\mathbf{x}_s, \mathbf{t}_l) \gamma_{l_s} + \sum_{l=1}^{d_1} \sum_{s=0}^{d_2} p_{\lambda_n} \left(\|\gamma_{l_s}\|_{\mathbf{W}_{l_s}} \right).$$

Now, let $\mathbf{Z}_{l_s}^* = \mathbf{Z}_{l_s} \mathbf{W}_{l_s}^{-\frac{1}{2}}$ and $\gamma_{l_s}^* = \mathbf{W}_{l_s} \gamma_{l_s}$. Consequently one has $\mathbf{W}_{l_s}^* = \frac{\mathbf{Z}_{l_s}^{*T} \mathbf{Z}_{l_s}^*}{n} = \mathbf{I}_{J_n}$ for $s = 0$, and $\mathbf{W}_{l_0}^* = 1$. Therefore, (3.3) can be rewritten as

$$\begin{aligned} \hat{\gamma}^* &= \underset{\gamma^*}{\operatorname{argmin}} L_n(\gamma^*) \\ &= \underset{\gamma^*}{\operatorname{argmin}} \left\{ \frac{1}{2n} \mathbf{Y} - \sum_{l=1}^{d_1} \sum_{s=0}^{d_2} \mathbf{Z}_{l_s}^*(\mathbf{x}_s, \mathbf{t}_l) \gamma_{l_s}^* + \sum_{l=1}^{d_1} \sum_{s=0}^{d_2} p_{\lambda_n} \left(\|\gamma_{l_s}^*\|_{\mathbf{W}_{l_s}^*} \right) \right\} \\ &= \underset{\gamma^*}{\operatorname{argmin}} \left\{ \frac{1}{2n} \mathbf{Y} - \sum_{l=1}^{d_1} \sum_{s=0}^{d_2} \mathbf{Z}_{l_s}^*(\mathbf{x}_s, \mathbf{t}_l) \gamma_{l_s}^* + \sum_{l=1}^{d_1} \sum_{s=1}^{d_2} p_{\lambda_n} (\|\gamma_{l_s}\|_2) + \sum_{l=1}^{d_1} p_{\lambda_n} (|\gamma_{l_0}^*|) \right\}. \end{aligned}$$

Therefore, for any given index pair $(l, s) \in S$, the partial derivative

$$\frac{\partial L_n(\gamma^*)}{\partial \gamma_{l_s}^*} = -\frac{1}{n} \mathbf{Z}_{l_s}^{*T} \mathbf{Y} - \sum_{l'=1}^{d_1} \sum_{s'=0}^{d_2} \mathbf{Z}_{l'_s'}^*(\mathbf{x}_{s'}, \mathbf{t}_{l'}) \gamma_{l'_s'}^* + \partial p_{\lambda_n} (\|\gamma_{l_s}^*\|_2),$$

where $\partial p_{\lambda_n}(\|\gamma_{ls}^*\|_2)$ is the subgradient of $p_{\lambda_n}(\|\gamma_{ls}^*\|_2)$.

We denote $\mathbf{C}_{ls}^*(\gamma^*) = -\frac{1}{n}\mathbf{Z}_{ls}^{*T} \left[\mathbf{Y} - \sum_{l'=1}^{d_1} \sum_{s'=0}^{d_2} \mathbf{Z}_{l's'}^*(\mathbf{x}_{s'}, \mathbf{t}_{l'}) \gamma_{l's'}^* \right]$. Then the KKT local optimality condition suggests that, any γ^* satisfying the following two conditions must be a local minimum of our penalized objective function,

- (i) $\mathbf{C}_{ls}^*(\gamma^*) = \mathbf{0}, \|\gamma_{ls}^*\|_2 > a\lambda_n$, for $(l, s) \in S^{(0)}$;
- (ii) $\|\mathbf{C}_{ls}^*(\gamma^*)\|_2 \leq \lambda_n, \|\gamma_{ls}^*\|_2 < \lambda_n$, for $(l, s) \notin S^{(0)}$.

Equivalently in terms of γ , the vector of untransformed coefficients, the sufficient conditions for a solution to be a local minimum are

$$(i) \quad \mathbf{C}_{ls}(\gamma) = \mathbf{0}, \|\gamma_{ls}\|_{\mathbf{W}_{ls}} > a\lambda_n, \text{ for } (l, s) \in S^{(0)}; \quad (3.10)$$

$$(ii) \quad \|\mathbf{C}_{ls}(\gamma)\|_{\mathbf{W}_{ls}^{-1}} \leq \lambda_n, \|\gamma_{ls}\|_{\mathbf{W}_{ls}} < \lambda_n, \text{ for } (l, s) \notin S^{(0)}. \quad (3.11)$$

Therefore, to prove $\hat{\gamma}^{(0)} \in A_n(\lambda_n)$, one only needs to prove (3.10) and (3.11) for $\gamma = \hat{\gamma}^{(0)}$.

When $\gamma = \hat{\gamma}^{(0)}$, the first equation in (3.10) and the second inequality in (3.11) naturally hold by the definition of $\hat{\gamma}^{(0)}$. So one only needs to prove that, except on an event whose probability tends to 0, as $n \rightarrow \infty$,

$$(i) \quad \hat{\gamma}_{ls}^{(0)} \mathbf{W}_{ls} > a\lambda_n, \text{ for } (l, s) \in S^{(0)}; \quad (3.12)$$

$$(ii) \quad \mathbf{C}_{ls}(\hat{\gamma}_{ls}^{(0)}) \mathbf{W}_{ls}^{-1} \leq \lambda_n, \text{ for } (l, s) \notin S^{(0)}. \quad (3.13)$$

We first prove (3.12). Note that $\hat{\gamma}_{ls}^{(0)} \mathbf{W}_{ls} = \hat{\alpha}_{ls}^{(0)}(X_s) T_l \frac{2}{n}$ is the empirical norm of a non-zero additive term of the oracle estimator. Lemma 3.7.3 implies that,

$$\begin{aligned} \hat{\gamma}_{ls}^{(0)} \mathbf{W}_{ls} &= \hat{\alpha}_{ls}^{(0)}(X_s) T_l \frac{\hat{\alpha}_{ls}^{(0)}(X_s) T_l}{\hat{\alpha}_{ls}^{(0)}(X_s) T_l} \frac{n}{n} \\ &\geq [\|\alpha_{ls}(X_s) T_l\|_2 - o_P(1)] \frac{\hat{\alpha}_{ls}^{(0)}(X_s) T_l}{\hat{\alpha}_{ls}^{(0)}(X_s) T_l}. \end{aligned}$$

Furthermore, Assumptions (C.2), (C.6) and Lemma A.4 in Xue and Yang (2006b) imply that the second multiplicative term converges to 1 in probability. Together with Assumptions (C.1), (C.3), and (iii) of Assumption (C.5), one has that, with probability goes to 1,

$$\begin{aligned}
\widehat{\gamma}_{ls}^{(0)} \quad W_{ls} &\geq \|\alpha_{ls}(X_s)T_l\|_2 - o_P(1) \\
&= \sqrt{E[\alpha_{ls}^2(X_s)T_l^2]} - o_P(1) \\
&= \sqrt{E[\alpha_{ls}^2(X_s)E\{T_l^2|X_s\}]} - o_P(1) \\
&\geq c\sqrt{E[\alpha_{ls}^2(X_s)]} - o_P(1) \\
&= c - o_P(1) \geq a\lambda_n.
\end{aligned} \tag{3.14}$$

Therefore (3.12) is proved. Now for (3.13), we define $\mathbf{Z}_{(1)} = (\mathbf{Z}_{ls}, (l, s) \in S^{(0)})$ as the column-wise combination of all \mathbf{Z}_{ls} matrices corresponding to all “nonzero” components, and $\mathbf{Z}_{(2)} = (\mathbf{Z}_{ls}, (l, s) \notin S^{(0)})$ as the column-wise combination of all \mathbf{Z}_{ls} matrices corresponding to all “redundant” components. Recall that $Y = \sum_{l=1}^{d_1} \sum_{s \in S_l} \alpha_{ls}(X_s)T_l + \varepsilon$. One has

$$\begin{aligned}
\mathbf{C}_{ls}(\widehat{\gamma}^{(0)}) &= \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{I}_n - \mathbf{Z}_{(1)} \left(\mathbf{Z}_{(1)}^T \mathbf{Z}_{(1)} \right)^{-1} \mathbf{Z}_{(1)}^T \mathbf{Y} \\
&= \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n [\delta + \varepsilon],
\end{aligned}$$

where $\mathbf{H}_n = \mathbf{I}_n - \mathbf{Z}_{(1)} \left(\mathbf{Z}_{(1)}^T \mathbf{Z}_{(1)} \right)^{-1} \mathbf{Z}_{(1)}^T$ and $\delta = \sum_{(l,s) \in S^{(0)}} \delta_{ls}$ with $\delta_{ls} = (\delta_{1,ls}, \dots, \delta_{n,ls})^T$ and $\delta_{i,ls} = \alpha_{ls}(x_{i,s})t_{i,l} - \widehat{\alpha}_{ls}^{(0)}(x_{i,s})t_{i,l}$. Therefore,

$$\begin{aligned}
P \left(\mathbf{C}_{ls}(\widehat{\gamma}^{(0)}) \quad \mathbf{w}_{ls}^{-1} > \lambda_n, \exists (l, s) \notin S^{(0)} \right) &\leq P \left(\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \delta \quad \mathbf{w}_{ls}^{-1} > \frac{\lambda_n}{2} \right) \\
&\quad + P \left(\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \varepsilon \quad \mathbf{w}_{ls}^{-1} > \frac{\lambda_n}{2} \right).
\end{aligned} \tag{3.15}$$

According to Lemma 3.7.1, one has

$$\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \delta \quad \mathbf{w}_{ls}^{-1} \leq \max_{(l,s) \notin S^{(0)}} \frac{c}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \delta \quad \leq \frac{c}{\sqrt{n}} \|\mathbf{H}_n \delta\|_2.$$

Note that $\mathbf{H}_n \leq \mathbf{I}_n$. That is, $\mathbf{I}_n - \mathbf{H}_n$ is semi-positive definite. The approximation theory in de Boor (2001) (p.149) gives that $\|\delta\|_2 \leq c\sqrt{n}N_n^{-(p+1)}$. Therefore, one has

$$\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \delta \mathbf{w}_{ls}^{-1} \leq \frac{1}{\sqrt{nc_1}} \|\delta\|_2 \leq cN_n^{-(p+1)}.$$

Consequently, (i) of Assumption (C.5) entails that

$$\lim_{n \rightarrow \infty} P \left(\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \delta \mathbf{w}_{ls}^{-1} > \frac{\lambda_n}{2} \right) = 0. \quad (3.16)$$

Similarly, one has $\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \varepsilon \mathbf{w}_{ls}^{-1} \leq \frac{1}{c_1 n} \max_{(l,s) \notin S^{(0)}} \mathbf{Z}_{ls}^T \mathbf{H}_n \varepsilon \mathbf{w}_{ls}^{-1}$. On the other hand, denote $\mathbf{H}^{(2)T} = \mathbf{Z}_{(2)}^T - \mathbf{Z}_{(2)}^T \mathbf{Z}_{(1)} \left(\mathbf{Z}_{(1)}^T \mathbf{Z}_{(1)} \right)^{-1} \mathbf{Z}_{(1)}^T$ and let $\mathbf{H}_{ls}^{(2)}$ be the J_n columns of $\mathbf{H}^{(2)}$ corresponding to \mathbf{Z}_{ls} with $\mathbf{H}_{ls}^{(2)} = \mathbf{Z}_{ls}^T \mathbf{H}_n$ for $(l,s) \in S^{(0)}$. Note that $\mathbf{H}^{(2)} \mathbf{H}^{(2)T} = \mathbf{Z}_{(2)} \mathbf{H}_n \mathbf{Z}_{(2)}^T \leq \mathbf{Z}_{(2)} \mathbf{Z}_{(2)}^T$. Therefore, one has $\mathbf{Z}_{ls}^T \mathbf{H}_n \varepsilon \mathbf{w}_{ls}^{-1} \leq \mathbf{Z}_{ls}^T \varepsilon \mathbf{w}_{ls}^{-1}$. Consequently,

$$\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \varepsilon \mathbf{w}_{ls}^{-1} \leq \frac{1}{c_1 n} \max_{(l,s) \notin S^{(0)}} \mathbf{Z}_{ls}^T \varepsilon \mathbf{w}_{ls}^{-1}.$$

According to Lemma 7 in Xue and Qu (2012),

$$E \left(\max_{(l,s) \notin S^{(0)}} \frac{1}{\sqrt{n}} \mathbf{Z}_{ls}^T \varepsilon \mathbf{w}_{ls}^{-1} \right) = O \left(\sqrt{\log(N_n d_1 d_2)} \right).$$

Then Markov inequality implies

$$P \left(\max_{(l,s) \notin S^{(0)}} \frac{1}{n} \mathbf{Z}_{ls}^T \mathbf{H}_n \varepsilon \mathbf{w}_{ls}^{-1} > \frac{\lambda_n}{2} \right) \leq C \sqrt{\frac{\log(N_n)}{n \lambda_n^2}}. \quad (3.17)$$

Then (3.13) follows from (3.16), (3.17) and (ii) of Assumption (C.5).

3.7.3 Proof of Theorem 3.3.2

To prove $\hat{\gamma}^{(0)}$ converges to the global minimum in (3.3), it suffices to show that

$$P \left(L_n(\gamma) \geq L_n(\hat{\gamma}^{(0)}) \text{ for all } \gamma \in \mathbb{R}^{d_1 d_2 J_n} \right) \rightarrow 1 \text{ as } n \rightarrow \infty. \quad (3.18)$$

We first define several subset of the index pair set $S = \{(l,s) | l = 1, \dots, d_1, s = 1, \dots, d_2\}$. Let $\mathcal{P} = S^{(0)}$ and $\mathcal{N} = S/S^{(0)} = \{(l,s) | (l,s) \notin S^{(0)}\}$. For any given $\gamma \in \mathbb{R}^{d_1 d_2 J_n}$, let

$$\begin{aligned} \mathcal{P}^+ &= \left\{ (l,s) \in \mathcal{P} = S^{(0)}, \|\gamma_{ls}\|_{\mathbf{w}_{ls}} > a\lambda_n \right\}, \quad \mathcal{P}^- = \mathcal{P}/\mathcal{P}^+; \\ \mathcal{N}^+ &= \left\{ (l,s) \in \mathcal{N} = S/S^{(0)}, \|\gamma_{ls}\|_{\mathbf{w}_{ls}} > \lambda_n \right\}, \quad \mathcal{N}^- = \mathcal{N}/\mathcal{N}^+. \end{aligned}$$

Zou, H. and Li, R. (2008). *One-step sparse estimates in nonconcave penalized likelihood models. *Annals of Statistics*, **36**, 1509–1533.*

