

# Notes on 'Adaptive Non-Parametric Regression With the $K$ -NN Fused Lasso'

Alden Green

February 14, 2019

Let  $\mathbf{X} = x_1, \dots, x_n$  be sampled i.i.d from  $\mu$  with density function  $p(\cdot)$  over some subset  $\mathcal{X}$  of Euclidean space, and suppose

$$y_i = f_0(x_i) + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d}{\sim} SG(\sigma^2)$$

holds for some unknown  $f_0$ . Let  $\hat{\theta}$  be the solution to the *fused lasso*

$$\hat{\theta} := \operatorname{argmin}_{\theta \in \mathbb{R}^n} \left\{ \frac{1}{2} \|y - \theta\|_2^2 + \lambda \|\nabla_G \theta\|_1 \right\}$$

where  $\lambda > 0$  is a tuning parameter, and  $\nabla_G$  is an oriented incidence matrix of the graph  $G$ .

The  $K$ -NN-FL estimator computes the fused lasso over the  $K$ -NN graph  $G_K$  of  $\mathbf{X}$ . The  $\epsilon$ -FL estimator computes the fused lasso over the  $\epsilon$  graph  $G_\epsilon$ .

The assumptions required for Theorems 1 and 2 are as follows.

(a) For all  $x \in \mathcal{X}$

$$0 < p_{\min} \leq p(x) \leq p_{\max} < \infty$$

(b) The base measure  $\mu$  in  $\mathcal{X}$  satisfies

$$r^d c_{1,d} \leq \mu(B_r(x)) \leq c_{2,d} r^d \quad (\forall x \in \mathcal{X})$$

(c) There exists a homeomorphism (continuous bijection with continuous inverse)  $h : \mathcal{X} \rightarrow [0, 1]^d$  such that

$$L_{\min} d_{\mathcal{X}}(x, x') \leq \|h(x) - h(x')\|_2 \leq L_{\max} d_{\mathcal{X}}(x, x') \quad (\forall x, x' \in \mathcal{X})$$

(d)  $g_0$  is piecewise Lipschitz<sup>1</sup>, meaning there exists a set  $\mathcal{S} \subset (0, 1)^d$  such that

$$(a) \quad \nu(\mathcal{S}) = 0.$$

---

<sup>1</sup>Technically, the requirement is slightly weaker than piecewise Lipschitz.

$$(b) \quad \mu\left(h^{-1}\left(S_\epsilon \cup ([0, 1]^d \setminus \Omega_\epsilon)\right)\right) \leq C_S \epsilon$$

(c) There exists a positive constant  $L_0$  such that if  $z$  and  $z'$  belong to the same connected component of  $\Omega_\epsilon \setminus B_\epsilon(\mathcal{S})$ , then

$$|g(z) - g(z')| \leq L_0 \|z - z'\|_2$$

where  $\Omega_\epsilon = [0, 1]^d \setminus B_\epsilon(\partial[0, 1]^d)$ .

**Theorem 1.** *Let  $K \asymp \log^{1+2r} n$  for some  $r > 0$ , Then under Assumptions 1-3, with an appropriate choice of the tuning parameter  $\lambda$ , the **K-NN-FL** estimator  $\hat{\theta}$  satisfies*

$$\left\| \hat{\theta} - \theta^\star \right\|_n^2 = O_{\mathbb{P}} \left( \frac{\log^{1+2r} n}{n} + \frac{\log^{1.5+r} n}{n} \|\nabla_{G_K} \theta^\star\|_1 \right)$$

*This upper bound also holds for  $\epsilon$ -NN-FL if we replace  $\|\nabla_{G_K} \theta^\star\|_1$  with  $\|\nabla_{G_\epsilon} \theta^\star\|_1$ .*

**Theorem 2.** *Under Assumptions 1-5, with an appropriate choice of the tuning parameter  $\lambda$ , the **K-NN-FL** estimator  $\hat{\theta}$  satisfies*

$$\left\| \hat{\theta} - \theta^\star \right\|_n^2 = \tilde{O}_{\mathbb{P}} \left( \frac{1}{n^{1/d}} \right).$$

## 1 Proofs

To ease proofs, we will assume  $\mathcal{X} = [0, 1]^d$ .

Construct  $G_{lat} = (V_{lat}, E_{lat})$  a lattice graph with equal side lengths in  $[0, 1]^d$ , where

$$V_{lat} = P_{lat}(N) := \left\{ \left( \frac{i_1}{N} - \frac{1}{2N}, \dots, \frac{i_d}{N} - \frac{1}{2N} \right) : i_1, \dots, i_d \in \{1, \dots, N\} \right\}$$

$$(z, z') \in E_{lat} \text{ if and only if } \|z - z'\| \leq \frac{1}{N}$$

where  $z$  and  $z' \in P_{lat}(N)$ .

Denoting  $I = P_{lat}$ , we define

$$P_I(x) = \operatorname{argmin} \{ \|x - z'\|_\infty, z' \in P_{lat}(N) \}$$

Then, let  $C(z) = \{x \in [0, 1]^d : z = P_I(x)\}$  be the collection of cells associated with the mesh  $P_{lat}(N)$ , noting that  $\{C(z) : z \in P_{lat}(N)\}$  defines a partition over  $[0, 1]^d$ .

**Quantization.** For a given  $\theta \in \mathbb{R}^n$ , the *quantization*  $\theta_I \in \mathbb{R}^n$

$$(\theta_I)_i := \theta_j, \quad \text{where } x_j = \underset{x_l, l \in [n]}{\operatorname{argmin}} \|P_I(x_i) - x_l\|_\infty$$

is constant over every cell  $C(z)$ . We now induce a signal in  $\mathbb{R}^{N^d}$  corresponding to the elements in  $I$ . Let  $\{z_1, \dots, z_{N^d}\} = I$ . Then we write

$$I_k = \{i \in [n] : P_I(x_i) = z_k\}$$

for  $k = 1, \dots, N^d$ . Define  $\theta^I \in \mathbb{R}^{N^d}$  by

$$(\theta^I)_k := \begin{cases} (\theta_I)_i, & x_i \in I_k \\ 0, & I_k = \emptyset \end{cases}$$

where we note that  $(\theta^I)$  is well-defined since  $(\theta_I)_i = (\theta_I)_j$  if  $x_i$  and  $x_j$  are both in  $I_k$ .

## 1.1 Controlling counts of mesh

Define the event  $\Omega$  as: “If  $x_i \in C(z_k)$  and  $x_i \in C(z_l)$  for  $z_k, z_l \in I$  with  $\|z_k - z_l\|_2 \leq \frac{1}{N}$ , then  $x_i$  and  $x_j$  are connected in the  $K$ -NN graph.” Then,

**Lemma 1.** *Take Assumptions 1-3, and additionally assume that  $N$  in the construction of  $G_{lat}(N)$  is chosen as*

$$N \geq \left\lceil \frac{3\sqrt{d}(2c_{2,d}p_{\max})^{1/d}n^{1/d}}{L_{\min}K^{1/d}} \right\rceil. \quad (1)$$

Then,

$$\mathbb{P}(\Omega) \geq 1 - n \exp(-K/3).$$

## 1.2 Bounding Empirical Process

**Lemma 2.**

## 1.3 Mesh embedding for $K$ -NN graph

**Lemma 3.** *Fix  $N$  to satisfy (1), and let us assume that the event  $\Omega$  from Lemma 1 holds. Denote  $I = P_{lat}(N)$  to be the mesh. Then, for all  $e \in \mathbb{R}^n$ , it holds that*

$$|e^T(\theta - \theta_I)| \leq 2\|e\|_\infty \|\nabla_{G_K}\theta\|_1, \quad (\forall \theta \in \mathbb{R}^n)$$

Moreover,

$$\|D\theta^I\|_1 \leq \|\nabla_{G_K}\theta\|_1, \quad (\forall \theta \in \mathbb{R}^n)$$

where  $D$  is the incidence matrix of  $G_{lat}$ .

*Proof.* Clearly

$$\langle \epsilon^T, \theta - \theta_I \rangle \leq \|\epsilon\|_\infty \cdot \|\theta - \theta_I\|_1$$

Then, for every  $i = 1, \dots, n$ , the event  $\Omega$  implies that there exists a  $j \in [n]$  such that

$$(\theta_I)_i = \theta_j, \quad (x_i, x_j) \in E_{G_K}$$

and therefore

$$\|\theta - \theta_I\|_1 \leq \|\nabla_{G_K} \theta\|_1$$

□

## 1.4 Bounding empirical process

**Lemma 4.** *Conditional on the event  $\Omega$ , we have that*

$$\langle \epsilon, \hat{\theta}_I - \hat{\theta}_I^* \rangle \leq \max_{u \in I} \sqrt{|C(u)|} \left( \|\Pi \tilde{\epsilon}\|_2 \left\| \hat{\theta} - \theta^* \right\|_2 + \|(D^\dagger)^T \tilde{\epsilon}\|_\infty \left[ \left\| \nabla_{G_K} \hat{\theta} \right\|_1 + \left\| \nabla_{G_K} \theta^* \right\|_1 \right] \right)$$

where  $\tilde{\epsilon}$  is an independent, mean-zero vector of subgaussian random variables,  $|C(u)| := \sum_{i \in [n]} \mathbb{I}(x_i \in C(u))$ , and  $D^\dagger$  is the pseudoinverse of the incidence matrix  $D$  of  $G_{\text{lat}}$ .

*Proof.* Writing

$$\tilde{\epsilon}_l = \left[ \max_{u \in I} |C(u)| \right]^{-1/2} \sum_{x_j \in I_l} \epsilon_j$$

we have

$$\langle \epsilon, \hat{\theta}_I - \theta_I^* \rangle = \left[ \max_{u \in I} |C(u)| \right]^{1/2} \langle \tilde{\epsilon}, \hat{\theta}^I - \theta^{*,I} \rangle$$

Now, divide up  $\tilde{\epsilon}$  into

$$\tilde{\epsilon} = P_1(\tilde{\epsilon}) + P_{1^\perp}(\tilde{\epsilon})$$

where  $P_1$  is the projection onto the span of  $\mathbf{1}$  the constant vector, and  $P_{1^\perp}$  the projection onto the space orthogonal to  $\mathbf{1}$ . Note that  $P_{1^\perp}(x) = (D^\dagger D)^T x$ .

Then, we have

$$\begin{aligned} \langle P_{1^\perp}(\tilde{\epsilon}), \hat{\theta}^I - \theta^{*,I} \rangle &= \langle (D^\dagger D)^T \tilde{\epsilon}, \hat{\theta}^I - \theta^{*,I} \rangle \\ &= \langle (D^\dagger)^T \tilde{\epsilon}, D(\hat{\theta}^I - \theta^{*,I}) \rangle \\ &\leq \|(D^\dagger)^T \tilde{\epsilon}\|_\infty \left\| D(\hat{\theta}^I - \theta^{*,I}) \right\|_1 \leq \|(D^\dagger)^T \tilde{\epsilon}\|_\infty \left[ \left\| \nabla_{G_K} \hat{\theta} \right\|_1 + \left\| \nabla_{G_K} \theta^* \right\|_1 \right] \end{aligned}$$

where the last inequality follows from the triangle inequality and Lemma 3.

On the other hand,

$$\langle P_1(\tilde{\epsilon}), \hat{\theta}^I - \theta^{*,I} \rangle \leq \|P_1(\tilde{\epsilon})\|_2 \left\| \hat{\theta}^I - \theta^{*,I} \right\|_2$$

and so the desired result follows. □

## 1.5 Proof of Theorem 1

We begin with a basic inequality

$$\frac{1}{2} \left\| \hat{\theta} - \theta^* \right\|_n^2 \leq \frac{1}{n} \langle \epsilon, \hat{\theta} - \theta^* \rangle + \lambda_n \left( \left\| \nabla_G \theta^* \right\|_1 - \left\| \nabla_G \hat{\theta} \right\|_1 \right)$$

We split up the **empirical process**,

$$\langle \epsilon, \hat{\theta} - \theta^* \rangle = \langle \epsilon, \hat{\theta} - \hat{\theta}_I \rangle + \langle \epsilon, \hat{\theta}_I - \theta_I^* \rangle + \langle \epsilon, \theta_I^* - \theta^* \rangle$$

Hereafter in the proof, we condition on the event  $\Omega$ . Lemma 3 gives us bounds on the first and third terms

$$\begin{aligned} \langle \epsilon, \hat{\theta} - \hat{\theta}_I \rangle &\leq 2 \|\epsilon\|_\infty \cdot \left\| \nabla_{G_K} \hat{\theta} \right\|_1 \\ \langle \epsilon, \theta_I^* - \theta^* \rangle &\leq 2 \|\epsilon\|_\infty \cdot \left\| \nabla_{G_K} \theta^* \right\|_1 \end{aligned}$$