

# 1 Appendix

## 1.1 Gradient Descent

---

**Algorithm 1** (WEIRDLY NUMBERED) Updated Algorithm 1

---

Initialize  $\boldsymbol{\lambda}^{(0)}$ .

**for** each iteration  $k = 0, 1, \dots$  until stopping criteria is reached **do**

Solve for  $\hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}^{(k)}) = \arg \min_{\boldsymbol{\theta} \in \Theta} L_T(\boldsymbol{\theta}, \boldsymbol{\lambda}^{(k)})$ .

Calculate the derivative of the model parameters with respect to the regularization parameters

$$\frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) = - \left[ \left[ \nabla_{\boldsymbol{\theta}}^2 \left( L(\mathbf{y}_T, f_{\boldsymbol{\theta}}(\mathbf{X}_T)) + \sum_{i=1}^J \lambda_i P_i(\boldsymbol{\theta}) \right) \right]^{-1} \nabla_{\boldsymbol{\theta}} P(\boldsymbol{\theta}) \right] \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}^{(k)})} \quad (1)$$

Calculate the gradient

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{y}_V, f_{\hat{\boldsymbol{\theta}}(\boldsymbol{\lambda})}(\mathbf{X}_V)) \Big|_{\boldsymbol{\lambda}=\boldsymbol{\lambda}^{(k)}} = \left[ \frac{\partial}{\partial \boldsymbol{\theta}} L(\mathbf{y}_V, f_{\boldsymbol{\theta}}(\mathbf{X}_V)) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}^{(k)})} \right]^{\top} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \Big|_{\boldsymbol{\lambda}=\boldsymbol{\lambda}^{(k)}} \quad (2)$$

Perform gradient step with step size  $t^{(k)}$

$$\boldsymbol{\lambda}^{(k+1)} := \boldsymbol{\lambda}^{(k)} - t^{(k)} \nabla_{\boldsymbol{\lambda}} L(\mathbf{y}_V, f_{\hat{\boldsymbol{\theta}}(\boldsymbol{\lambda})}(\mathbf{X}_V)) \Big|_{\boldsymbol{\lambda}=\boldsymbol{\lambda}^{(k)}} \quad (3)$$

**end for**

---

### 1.1.1 $K$ -fold Cross Validation

We can perform joint optimization for  $K$ -fold cross validation by reformulating the problem. Let  $(\mathbf{y}, \mathbf{X})$  be the full data set. We denote the  $k$ th fold as  $(\mathbf{y}_k, \mathbf{X}_k)$  and its complement as  $(\mathbf{y}_{-k}, \mathbf{X}_{-k})$ . Then the objective of this joint optimization problem is the average validation cost across all  $K$  folds:

$$\begin{aligned} & \arg \min_{\boldsymbol{\lambda} \in \Lambda} \frac{1}{K} \sum_{k=1}^K L(\mathbf{y}_k, f_{\hat{\boldsymbol{\theta}}^{(k)}(\boldsymbol{\lambda})}(\mathbf{X}_k)) \\ \text{s.t. } & \hat{\boldsymbol{\theta}}^{(k)}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\theta} \in \Theta} L(\mathbf{y}_{-k}, f_{\boldsymbol{\theta}}(\mathbf{X}_{-k})) + \sum_{i=1}^J \lambda_i P_i(\boldsymbol{\theta}) \text{ for } k = 1, \dots, K \end{aligned} \quad (4)$$

## 1.2 Proof of Theorem 1

*Proof.* We will show that for a given  $\boldsymbol{\lambda}_0$  that satisfies the given conditions, the validation loss is continuously differentiable within some neighborhood of  $\boldsymbol{\lambda}_0$ . It then follows that if the theorem conditions hold true for almost every  $\boldsymbol{\lambda}$ , then the validation loss is continuously differentiable with respect to  $\boldsymbol{\lambda}$  at almost every  $\boldsymbol{\lambda}$ .

Suppose the theorem conditions are satisfied at  $\lambda_0$ . Let  $\mathbf{B}'$  be an orthonormal set of basis vectors that span the differentiable space  $\Omega^{L_T}(\hat{\boldsymbol{\theta}}(\lambda_0), \lambda_0)$  with the subset of vectors  $\mathbf{B}$  that span the model parameter space.

Let  $\tilde{L}_T(\boldsymbol{\theta}, \lambda)$  be the gradient of  $L_T(\cdot, \lambda)$  at  $\boldsymbol{\theta}$  with respect to the basis  $\mathbf{B}$ :

$$\tilde{L}_T(\boldsymbol{\theta}, \lambda) = {}_{\mathbf{B}} \nabla L_T(\cdot, \lambda)|_{\boldsymbol{\theta}} \quad (5)$$

Since  $\hat{\boldsymbol{\theta}}(\lambda_0)$  is the minimizer of the training loss, the gradient of  $L_T(\cdot, \lambda_0)$  with respect to the basis  $\mathbf{B}$  must be zero at  $\hat{\boldsymbol{\theta}}(\lambda_0)$ :

$${}_{\mathbf{B}} \nabla L_T(\cdot, \lambda_0)|_{\hat{\boldsymbol{\theta}}(\lambda_0)} = \tilde{L}_T(\hat{\boldsymbol{\theta}}(\lambda_0), \lambda_0) = 0 \quad (6)$$

From our assumptions, we know that there exists a neighborhood  $W$  containing  $\lambda_0$  such that  $\tilde{L}_T$  is continuously differentiable along directions in the differentiable space  $\Omega^{L_T}(\hat{\boldsymbol{\theta}}(\lambda_0), \lambda_0)$ . Also, the Jacobian matrix  $D\tilde{L}_T(\cdot, \lambda_0)|_{\hat{\boldsymbol{\theta}}(\lambda_0)}$  with respect to basis  $\mathbf{B}$  is nonsingular. Therefore, by the implicit function theorem, there exist open sets  $U \subseteq W$  containing  $\lambda_0$  and  $V$  containing  $\hat{\boldsymbol{\theta}}(\lambda_0)$  and a continuously differentiable function  $\gamma : U \rightarrow V$  such that for every  $\lambda \in U$ , we have that

$$\tilde{L}_T(\gamma(\lambda), \lambda) = \nabla_{\mathbf{B}} L_T(\cdot, \lambda)|_{\gamma(\lambda)} = 0 \quad (7)$$

That is, we know that  $\gamma(\lambda)$  is a continuously differentiable function that minimizes  $L_T(\cdot, \lambda)$  in the differentiable space  $\Omega^{L_T}(\hat{\boldsymbol{\theta}}(\lambda_0), \lambda_0)$ . Since we assumed that the differentiable space is a local optimality space of  $L_T(\cdot, \lambda)$  in the neighborhood  $W$ , then for every  $\lambda \in U$ ,

$$\hat{\boldsymbol{\theta}}(\lambda) = \arg \min_{\boldsymbol{\theta}} L_T(\boldsymbol{\theta}, \lambda) = \arg \min_{\boldsymbol{\theta} \in \Omega^{L_T}(\hat{\boldsymbol{\theta}}(\lambda_0), \lambda_0)} L_T(\boldsymbol{\theta}, \lambda) = \gamma(\lambda) \quad (8)$$

Therefore, we have shown that if  $\lambda_0$  satisfies the assumptions given in the theorem, the fitted model parameters  $\hat{\boldsymbol{\theta}}(\lambda)$  is a continuously differentiable function within a neighborhood of  $\lambda_0$ . We can then apply the chain rule to get the gradient of the validation loss.  $\square$

## 1.3 Examples - Checking Theorem 1 Conditions

### 1.3.1 Elastic Net

Next, we show that the joint optimization problem satisfies all three conditions in Theorem 1:

Condition 1: The nonzero indices of the elastic net estimates stay locally constant for almost every  $\lambda$ . Therefore,  $S_\lambda$  is a local optimality space for  $L_T(\cdot, \lambda)$   $\checkmark$

Condition 2: The  $\ell_1$  penalty is smooth when restricted to  $S_\lambda$ .  $\checkmark$

Condition 3: The Hessian matrix of  $L_T(\cdot, \lambda)$  with respect to the columns of  $\mathbf{I}_{I(\lambda)}$  is  $\mathbf{I}_{I(\lambda)}^\top \mathbf{X}_T^\top \mathbf{X}_T \mathbf{I}_{I(\lambda)} + \lambda_2 \mathbf{I}$ . This is positive definite if  $\lambda_2 > 0$ .  $\checkmark$

### 1.3.2 Sparse Group Lasso

Since the reasoning for the first two conditions is exactly the same as in the Elastic Net (Section 1.3.1), we just give the calculations for the third condition.

Condition 3: The Hessian matrix of  $L_T(\cdot, \boldsymbol{\lambda})$  with respect to the columns of  $\mathbf{I}_{I(\boldsymbol{\lambda})}$  is

$$\frac{1}{n} \mathbf{I}_{I(\boldsymbol{\lambda})}^\top \mathbf{X}_T^\top \mathbf{X}_T \mathbf{I}_{I(\boldsymbol{\lambda})} + \lambda_1 \mathbf{B}(\boldsymbol{\lambda}) + \epsilon \mathbf{I}_p \quad (9)$$

where  $\mathbf{B}(\boldsymbol{\lambda})$  is a block diagonal matrix with components

$$\left\| \tilde{\boldsymbol{\theta}}^{(m)}(\boldsymbol{\lambda}) \right\|_2^{-1} \left( \mathbf{I} - \frac{\tilde{\boldsymbol{\theta}}^{(m)}(\boldsymbol{\lambda}) \tilde{\boldsymbol{\theta}}^{(m)}(\boldsymbol{\lambda})^\top}{\left\| \tilde{\boldsymbol{\theta}}^{(m)}(\boldsymbol{\lambda}) \right\|_2^2} \right) \quad (10)$$

for  $m = 1, \dots, M$  from top left to bottom right. The Hessian is positive definite for any fixed  $\epsilon > 0$ . ✓

### 1.3.3 Generalized Lasso (condition check)

The first two conditions in Theorem 1 are satisfied by similar reasoning to that discussed in Section 1.3.1. For the third condition, we need to check that the Hessian matrix is invertible.

Condition 3: The Hessian matrix of  $L_T(\cdot, \boldsymbol{\lambda})$  with respect to  $\mathbf{U}_\lambda$  is

$$\mathbf{U}_\lambda^\top \mathbf{X}_T^\top \mathbf{X}_T \mathbf{U}_\lambda + \epsilon \mathbf{U}_\lambda \quad (11)$$

This is positive definite for any fixed  $\epsilon > 0$ . ✓

### 1.3.4 Generalized Lasso (the whole thing)

The generalized lasso (?) penalizes the  $\ell_1$  norm of the coefficients  $\boldsymbol{\theta}$  weighted by some matrix  $\mathbf{D}$ . Depending on the choice of  $\mathbf{D}$ , the generalized lasso induces different structural constraints on the regression coefficients. Special cases include the fused lasso, trend filtering, and wavelet smoothing (?), (?), (?).

To tune the regularization parameter  $\lambda$ , we formulate the generalized lasso as a joint optimization problem:

$$\begin{aligned} & \min_{\lambda \in \mathbb{R}_+} \frac{1}{2} \|\mathbf{y}_V - \mathbf{X}_V \hat{\boldsymbol{\theta}}(\lambda)\|^2 \\ \text{s.t. } & \hat{\boldsymbol{\theta}}(\lambda) = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y}_T - \mathbf{X}_T \boldsymbol{\theta}\|^2 + \lambda \|\mathbf{D} \boldsymbol{\theta}\|_1 + \frac{1}{2} \epsilon \|\boldsymbol{\theta}\|_2^2 \end{aligned} \quad (12)$$

Let  $I(\lambda)$  denote the indices of the zero elements of  $\mathbf{D} \hat{\boldsymbol{\theta}}(\lambda)$ :

$$I(\lambda) = \left\{ i \mid (\mathbf{D} \hat{\boldsymbol{\theta}}(\lambda))_i = 0 \text{ for } i = 1, \dots, p \right\} \quad (13)$$

Since  $\|\mathbf{D} \boldsymbol{\theta}\|_1$  is differentiable in  $\boldsymbol{\theta}$  only along directions where the current zero elements of  $\mathbf{D} \boldsymbol{\theta}$  remain zero, the differentiable space  $S_\lambda$  is the null space of  $\mathbf{I}_{I(\lambda)}^\top \mathbf{D}$ , denoted  $\mathcal{N}(\mathbf{I}_{I(\lambda)}^\top \mathbf{D})$ . Let  $\mathbf{U}_\lambda$  be an orthonormal basis for  $\mathcal{N}(\mathbf{I}_{I(\lambda)}^\top \mathbf{D})$ .

Now we show the gradient calculations. Following Algorithm 2, we first define the locally equivalent joint optimization problem:

$$\begin{aligned} & \min_{\lambda \in \mathbb{R}_+} \frac{1}{2} \|\mathbf{y}_V - \mathbf{X}_V \mathbf{U}_\lambda \hat{\boldsymbol{\beta}}(\lambda)\|^2 \\ \text{s.t. } & \hat{\boldsymbol{\beta}}(\lambda) = \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y}_T - \mathbf{X}_T \mathbf{U}_\lambda \boldsymbol{\beta}\|^2 + \lambda \|\mathbf{D} \mathbf{U}_\lambda \boldsymbol{\beta}\|_1 + \frac{1}{2} \epsilon \|\mathbf{U}_\lambda \boldsymbol{\beta}\|_2^2 \end{aligned} \quad (14)$$

Implicit differentiation with respect to  $\lambda$  of the gradient condition for  $\hat{\boldsymbol{\beta}}(\lambda)$  gives us

$$\frac{\partial}{\partial \lambda} \hat{\boldsymbol{\beta}}(\lambda) = -(\mathbf{U}_\lambda^\top \mathbf{X}_T^\top \mathbf{X}_T \mathbf{U}_\lambda + \epsilon \mathbf{U}_\lambda)^\top \mathbf{U}_\lambda^\top \mathbf{D}^\top \text{sgn}(\mathbf{D} \mathbf{U}_\lambda \hat{\boldsymbol{\beta}}(\lambda)) \quad (15)$$

The chain rule then gives the gradient of the validation loss with respect to  $\lambda$ :

$$\nabla_\lambda L(\mathbf{y}_V, f_{\hat{\boldsymbol{\theta}}(\lambda)}(\mathbf{X}_V)) = - \left( \mathbf{X}_V \mathbf{U}_\lambda \frac{\partial}{\partial \lambda} \hat{\boldsymbol{\beta}}(\lambda) \right)^\top (\mathbf{y}_V - \mathbf{X}_V \mathbf{U}_\lambda \hat{\boldsymbol{\beta}}(\lambda)) \quad (16)$$

### 1.3.5 Ugly Math from Additive Models

To formalize our optimization problem we give a bit of notation. Let  $\mathbf{Z}_T \in \mathbb{R}^{|T| \times q}$ ,  $\mathbf{Z}_V \in \mathbb{R}^{|V| \times q}$  be the linear covariates from the training and validation sets. Let  $\mathbf{X} \in \mathbb{R}^{n \times p}$  be the design matrix for the nonlinear covariates where the first  $|T|$  observations are from the training set and the last  $|V|$  observations are from the validation set. Let  $\mathbf{I}_T$  and  $\mathbf{I}_V$  be matrices such that  $\mathbf{X}_T = \mathbf{I}_T \mathbf{X}$  and  $\mathbf{X}_V = \mathbf{I}_V \mathbf{X}$ . We combine the validation and training nonlinear covariates into one design matrix since we use the training data to fit estimates for  $f_i$  at the validation points.

The joint optimization problem is

$$\begin{aligned} & \min_{\boldsymbol{\lambda} \in \mathbb{R}_+^{p+1}} \frac{1}{2} \left\| \mathbf{y}_V - \mathbf{Z}_V \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) - \mathbf{I}_V \sum_{i=1}^p \hat{\boldsymbol{\theta}}^{(i)}(\boldsymbol{\lambda}) \right\|_2^2 \\ \text{s.t. } & \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}), \hat{\boldsymbol{\theta}}^{(i)}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\beta}, \boldsymbol{\theta}} \frac{1}{2} \left\| \mathbf{y}_T - \mathbf{Z}_T \boldsymbol{\beta} - \mathbf{I}_T \sum_{i=1}^p \boldsymbol{\theta}^{(i)} \right\|_2^2 + \lambda_0 \|\boldsymbol{\beta}\|_1 \\ & + \frac{1}{2} \sum_{i=1}^p \lambda_i \left\| \mathbf{D}_{\mathbf{x}_i}^{(2)} \boldsymbol{\theta}^{(i)} \right\|_2^2 + \frac{1}{2} \epsilon \left( \|\boldsymbol{\beta}\|_2^2 + \sum_{i=1}^p \|\boldsymbol{\theta}^{(i)}\|_2^2 \right) \end{aligned} \quad (17)$$

The matrix  $\mathbf{D}_{\mathbf{x}_i}^{(2)}$  gives the second-order differences between the nonparametric estimates of  $f_i$  for unevenly-spaced inputs. Construction of  $\mathbf{D}_{\mathbf{x}_i}^{(2)}$  is given in the Appendix. Note that we have presented an extremely general form on the training criterion, where there are separate  $\lambda_i$  for  $\left\| \mathbf{D}_{\mathbf{x}_i}^{(2)} \boldsymbol{\theta}^{(i)} \right\|_2^2$ . In practice, one would pool  $\lambda_i$  since tuning more than 2 penalty parameters is difficult.

In this example, the lasso is the only penalty which is not everywhere differentiable. Let the nonzero indices of  $\hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})$  be denoted  $I(\boldsymbol{\lambda}) = \{i | \hat{\beta}_i(\boldsymbol{\lambda}) \neq 0 \text{ for } i = 1, \dots, p\}$ . The differentiable space is then  $S_{\boldsymbol{\lambda}} = \mathbf{C}(\mathbf{I}_{I(\boldsymbol{\lambda})}) \oplus \mathbb{R}^{n \times p}$ .

We now calculate the gradient of the validation loss. Given  $I(\boldsymbol{\lambda})$ , the nonzero set of  $\hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})$ , the locally equivalent joint optimization problem as

$$\begin{aligned} & \min_{\boldsymbol{\lambda} \in \mathbb{R}_+^2} \frac{1}{2} \left\| \mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) - \mathbf{I}_V \sum_{i=1}^p \hat{\boldsymbol{\theta}}^{(i)}(\boldsymbol{\lambda}) \right\|_2^2 \\ \text{s.t. } & \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}), \hat{\boldsymbol{\theta}}^{(i)}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\eta}, \boldsymbol{\theta}} \frac{1}{2} \left\| \mathbf{y}_T - \mathbf{X}_{T, I(\boldsymbol{\lambda})} \boldsymbol{\eta} - \mathbf{I}_T \sum_{i=1}^p \boldsymbol{\theta}^{(i)} \right\|_2^2 \\ & + \lambda_0 \|\boldsymbol{\eta}\|_1 + \frac{1}{2} \sum_{i=1}^p \lambda_i \left\| \mathbf{D}(\mathbf{z}) \boldsymbol{\theta}^{(i)} \right\|_2^2 + \frac{1}{2} \epsilon \left( \|\boldsymbol{\eta}\|_2^2 + \sum_{i=1}^p \|\boldsymbol{\theta}^{(i)}\|_2^2 \right) \end{aligned} \quad (18)$$

We follow the same steps as before to calculate the gradient. The gradient is then

$$\nabla_{\lambda_j} L_V(\boldsymbol{\lambda}) = - \left( \mathbf{X}_{V, I(\boldsymbol{\lambda})} \frac{\partial}{\partial \lambda_j} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) + \mathbf{I}_V \sum_{i=1}^p \frac{\partial}{\partial \lambda_j} \hat{\boldsymbol{\theta}}^{(i)}(\boldsymbol{\lambda}) \right)^\top \left( \mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) - \mathbf{I}_V \sum_{i=1}^p \hat{\boldsymbol{\theta}}^{(i)}(\boldsymbol{\lambda}) \right)$$

where

$$\begin{bmatrix} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) \\ \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\theta}}^{(1)}(\boldsymbol{\lambda}) \\ \dots \\ \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\theta}}^{(p)}(\boldsymbol{\lambda}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \lambda_0} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) & \dots & \frac{\partial}{\partial \lambda_p} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) \\ \frac{\partial}{\partial \lambda_0} \hat{\boldsymbol{\theta}}^{(1)}(\boldsymbol{\lambda}) & \dots & \frac{\partial}{\partial \lambda_p} \hat{\boldsymbol{\theta}}^{(1)}(\boldsymbol{\lambda}) \\ \dots & \dots & \dots \\ \frac{\partial}{\partial \lambda_0} \hat{\boldsymbol{\theta}}^{(p)}(\boldsymbol{\lambda}) & \dots & \frac{\partial}{\partial \lambda_p} \hat{\boldsymbol{\theta}}^{(p)}(\boldsymbol{\lambda}) \end{bmatrix} = H^{-1} \begin{bmatrix} \text{sgn}(\hat{\boldsymbol{\eta}}(\boldsymbol{\lambda})) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & D_{x_1}^T D_{x_1} \hat{\boldsymbol{\theta}}^{(1)}(\boldsymbol{\lambda}) & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & D_{x_p}^T D_{x_p} \hat{\boldsymbol{\theta}}^{(p)}(\boldsymbol{\lambda}) \end{bmatrix} \quad (19)$$

The hideous matrix  $H$  is given in the appendix.

The second-order difference matrix construction: Let the  $n$  observations be ordered according to the  $i$ th covariate such that  $x_{ik_1} \leq x_{ik_2} \leq \dots \leq x_{ik_n}$ .  $\mathbf{D}_{\mathbf{x}_i}^{(1)} \in \mathbb{R}^{n \times n}$  is the corresponding first-order difference matrix; so row  $j = 1, \dots, n-1$  of  $\mathbf{D}_{\mathbf{x}_i}^{(1)}$  has a -1 in position  $k_j$ , 1 in position  $k_{j+1}$ , and 0 elsewhere and row  $n$  is all zeros. Then  $\mathbf{D}_{\mathbf{x}_i}^{(2)}$  is

$$\mathbf{D}_{\mathbf{x}_i}^{(1)} \cdot \text{diag} \left( \frac{1}{x_{ik_2} - x_{ik_1}}, \frac{1}{x_{ik_3} - x_{ik_2}}, \dots, \frac{1}{x_{ik_n} - x_{ik_{n-1}}}, 0 \right) \cdot \mathbf{D}_{\mathbf{x}_i}^{(1)} \quad (20)$$

### 1.3.6 Additive Partially Linear Models

By the same reasoning in Section 1.3.1, the first two conditions of Theorem 1 are satisfied. We now check for the third condition.

Condition 3: The Hessian matrix of  $L_T(\cdot, \boldsymbol{\lambda})$  with respect to the basis

$$\begin{bmatrix} \mathbf{I}_{I(\boldsymbol{\lambda})} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \quad (21)$$

is

$$H = \begin{bmatrix} \mathbf{I}_{I(\boldsymbol{\lambda})}^\top \mathbf{X}_T^\top \mathbf{X}_T \mathbf{I}_{I(\boldsymbol{\lambda})} + \epsilon \mathbf{I} & \mathbf{I}_{I(\boldsymbol{\lambda})}^\top \mathbf{X}_T^\top \mathbf{I}_T \\ \mathbf{I}_T^\top \mathbf{X}_T \mathbf{I}_{I(\boldsymbol{\lambda})} & \mathbf{I}_T^\top \mathbf{I}_T + \lambda_2 \mathbf{D}(\mathbf{z})^\top \mathbf{D}(\mathbf{z}) + \epsilon \mathbf{I} \end{bmatrix} \quad (22)$$

The Hessian matrix is invertible for any  $\lambda_2 > 0$  and any fixed  $\epsilon > 0$ .

### 1.3.7 Additive Models with Heterogeneous Smoothness

The problem is composed solely of quadratic functions, so the validation loss  $L(\mathbf{y}_v, \mathbf{I}_V \sum_{i=1}^p \hat{\boldsymbol{\theta}}^{(i)}(\boldsymbol{\lambda}))$  differentiable everywhere with respect to  $\boldsymbol{\lambda}$ . Choosing  $\epsilon > 0$  will satisfy the positive definite condition.

## 1.4 Gradient Derivations

### 1.4.1 Un-pooled Sparse Group Lasso

The joint optimization formulation of the un-pooled sparse group lasso is

$$\begin{aligned} & \min_{\boldsymbol{\lambda} \in \mathbb{R}_+^2} \frac{1}{2n} \left\| \mathbf{y}_V - \mathbf{X}_V \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \right\|_2^2 \\ \text{s.t. } & \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\theta}} \frac{1}{2n} \left\| \mathbf{y}_T - \mathbf{X}_T \boldsymbol{\theta} \right\|_2^2 + \sum_{m=1}^M \lambda_1^{(m)} \left\| \boldsymbol{\theta}^{(m)} \right\|_2 + \lambda_2 \left\| \boldsymbol{\theta} \right\|_1 + \frac{1}{2} \epsilon \left\| \boldsymbol{\theta} \right\|_2^2 \end{aligned} \quad (23)$$

Let  $I(\boldsymbol{\lambda}) = \{i | \hat{\theta}_i(\boldsymbol{\lambda}) \neq 0 \text{ for } i = 1, \dots, p\}$ . With similar reasoning in Section 2.4.3, the differentiable space for this problem is  $\text{span}(\mathbf{I}_{I(\boldsymbol{\lambda})})$ . All three conditions of Theorem 1 are satisfied. We note that the Hessian in this problem is

$$\frac{1}{n} \mathbf{X}_{T, I(\boldsymbol{\lambda})}^\top \mathbf{X}_{T, I(\boldsymbol{\lambda})} + \mathbf{B}(\boldsymbol{\lambda}) + \epsilon \mathbf{I} \quad (24)$$

where  $\mathbf{B}(\boldsymbol{\lambda})$  is the block diagonal matrix with components  $m = 1, 2, \dots, M$

$$\frac{\lambda_1^{(m)}}{\left\| \boldsymbol{\theta}^{(m)} \right\|_2} \left( \mathbf{I} - \frac{1}{\left\| \boldsymbol{\theta}^{(m)} \right\|_2^2} \boldsymbol{\theta}^{(m)} \boldsymbol{\theta}^{(m)\top} \right) \quad (25)$$

from top left to bottom right. This is positive definite for any  $\epsilon > 0$ .

To find the gradient, the locally equivalent joint optimization with a smooth training criterion is

$$\begin{aligned} & \min_{\boldsymbol{\lambda} \in \mathbb{R}_+^2} \frac{1}{2n} \left\| \mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \right\|_2^2 \\ \text{s.t. } & \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\beta}} \frac{1}{2n} \left\| \mathbf{y}_T - \mathbf{X}_{T, I(\boldsymbol{\lambda})} \boldsymbol{\beta} \right\|_2^2 + \sum_{m=1}^M \lambda_1^{(m)} \left\| \boldsymbol{\beta}^{(m)} \right\|_2 + \lambda_2 \left\| \boldsymbol{\beta} \right\|_1 + \frac{1}{2} \epsilon \left\| \boldsymbol{\beta} \right\|_2^2 \end{aligned} \quad (26)$$

Implicit differentiation of the gradient condition with respect to the regularization parameters gives us

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) &= \begin{bmatrix} \frac{\partial}{\partial \lambda_1^{(1)}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) & \dots & \frac{\partial}{\partial \lambda_1^{(M)}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) & \frac{\partial}{\partial \lambda_2} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \end{bmatrix} \\ &= - \left( \frac{1}{n} \mathbf{X}_{T, I(\boldsymbol{\lambda})}^\top \mathbf{X}_{T, I(\boldsymbol{\lambda})} + \mathbf{B}(\boldsymbol{\lambda}) + \epsilon \mathbf{I} \right)^{-1} \begin{bmatrix} \mathbf{C}(\hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) & \text{sgn}(\hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) \end{bmatrix} \end{aligned} \quad (27)$$

where  $\mathbf{C}(\hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}))$  has columns  $m = 1, 2, \dots, M$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ \frac{\hat{\boldsymbol{\beta}}^{(m)}(\boldsymbol{\lambda})}{\left\| \hat{\boldsymbol{\beta}}^{(m)}(\boldsymbol{\lambda}) \right\|_2} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (28)$$

By the chain rule, we get that the gradient of the validation error is

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{y}_V, \mathbf{X}_V \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) = \frac{1}{n} \left( \mathbf{X}_{V, I(\boldsymbol{\lambda})} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \right)^\top (\mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) \quad (29)$$

### 1.4.2 Additive Partially Linear Model with three penalties

The joint optimization formulation of the additive partially linear model with the elastic net penalty for the linear model  $\beta$  and the H-P filter for the nonparametric estimates  $\theta$  is

$$\begin{aligned} & \min_{\lambda \in \mathbb{R}_+^2} \frac{1}{2} \left\| \mathbf{y}_V - \mathbf{X}_V \hat{\beta}(\lambda) - (\mathbf{I} - \mathbf{I}_T) \hat{\theta}(\lambda) \right\|_2^2 \\ \text{s.t. } & \hat{\beta}(\lambda), \hat{\theta}(\lambda) = \arg \min_{\beta, \theta} \frac{1}{2} \left\| \mathbf{y}_T - \mathbf{X}_T \beta - \mathbf{I}_T \theta \right\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{1}{2} \lambda_2 \|\beta\|_2^2 + \frac{1}{2} \lambda_3 \|D(\mathbf{z}) \theta\|_2^2 + \frac{1}{2} \epsilon \|\theta\|_2^2 \end{aligned} \quad (30)$$

The differentiable space is exactly the same as that given in Section ?? . Also, all three conditions of Theorem 1 are satisfied. Note that the Hessian of the training criterion with respect to the basis in (18) is

$$H = \begin{bmatrix} \mathbf{I}_{I(\lambda)}^\top \mathbf{X}_T^\top \mathbf{X}_T \mathbf{I}_{I(\lambda)} + \lambda_2 \mathbf{I} & \mathbf{I}_{I(\lambda)}^\top \mathbf{X}_T^\top \mathbf{I}_T \\ \mathbf{I}_T^\top \mathbf{X}_T \mathbf{I}_{I(\lambda)} & \mathbf{I}_T^\top \mathbf{I}_T + \lambda_3 D(\mathbf{z})^\top D(\mathbf{z}) + \epsilon \mathbf{I} \end{bmatrix} \quad (31)$$

To find the gradient, we first consider the locally equivalent joint optimization problem with a smooth training criterion:

$$\begin{aligned} & \min_{\lambda \in \mathbb{R}_+^2} \frac{1}{2} \left\| \mathbf{y}_V - \mathbf{X}_{V,I(\lambda)} \hat{\eta}(\lambda) - (\mathbf{I} - \mathbf{I}_T) \hat{\theta}(\lambda) \right\|_2^2 \\ \text{s.t. } & \hat{\eta}(\lambda), \hat{\theta}(\lambda) = \arg \min_{\eta, \theta} \frac{1}{2} \left\| \mathbf{y}_T - \mathbf{X}_{T,I(\lambda)} \eta - \mathbf{I}_T \theta \right\|_2^2 + \lambda_1 \|\eta\|_1 + \frac{1}{2} \lambda_2 \|\eta\|_2^2 + \frac{1}{2} \lambda_3 \|D(\mathbf{z}) \theta\|_2^2 + \frac{1}{2} \epsilon \|\theta\|_2^2 \end{aligned} \quad (32)$$

After implicit differentiation of the gradient condition with respect to the regularization parameters, we get that

$$\begin{bmatrix} \frac{\partial}{\partial \lambda} \hat{\eta}(\lambda) \\ \frac{\partial}{\partial \lambda} \hat{\theta}(\lambda) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \lambda_1} \hat{\eta}(\lambda) & \frac{\partial}{\partial \lambda_3} \hat{\eta}(\lambda) & \frac{\partial}{\partial \lambda_3} \hat{\eta}(\lambda) \\ \frac{\partial}{\partial \lambda_1} \hat{\theta}(\lambda) & \frac{\partial}{\partial \lambda_2} \hat{\theta}(\lambda) & \frac{\partial}{\partial \lambda_3} \hat{\theta}(\lambda) \end{bmatrix} = -H^{-1} \begin{bmatrix} \text{sgn}(\hat{\eta}(\lambda)) & \hat{\eta}(\lambda) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & D(\mathbf{z})^\top D(\mathbf{z}) \hat{\theta}(\lambda) \end{bmatrix} \quad (33)$$

We then apply the chain rule to get the gradient direction of the validation loss with respect to  $\lambda$

$$\nabla_{\lambda} L_V(\lambda) = - \left( \mathbf{X}_{V,I(\lambda)} \frac{\partial}{\partial \lambda} \hat{\eta}(\lambda) + (\mathbf{I} - \mathbf{I}_T) \frac{\partial}{\partial \lambda} \hat{\theta}(\lambda) \right)^\top \left( \mathbf{y}_V - \mathbf{X}_{V,I(\lambda)} \hat{\eta}(\lambda) - (\mathbf{I} - \mathbf{I}_T) \hat{\theta}(\lambda) \right) \quad (34)$$

## 1.5 Backtracking Line Search

Let the criterion function be  $L : \mathbb{R}^n \rightarrow \mathbb{R}$ . Suppose that the descent algorithm is currently at point  $x$  with descent direction  $\Delta x$ . Backtracking line search uses a heuristic for finding a step size  $t \in (0, 1]$  such that the value of the criterion is minimized. The method depends on constants  $\alpha \in (0, 0.5)$  and  $\beta \in (0, 1)$ .

---

**Algorithm 2** Backtracking Line Search

---

Initialize  $t = 1$ .  
**while**  $L(\mathbf{x} + t\Delta\mathbf{x}) > L(\mathbf{x}) + \alpha t \nabla L(\mathbf{x})^T \Delta\mathbf{x}$  **do**  
    Update  $t := \beta t$   
**end while**

---



## 1.6 Joint Optimization with Accelerated Gradient Descent and Adaptive Restarts

---

**Algorithm 3** Joint Optimization with Accelerated Gradient Descent and Adaptive Restarts

---

Initialize  $\lambda^{(0)}$ .

**while** stopping criteria is not reached **do**

**for** each iteration  $k = 0, 1, \dots$  **do**

    Solve for  $\hat{\theta}(\lambda^{(k)}) = \arg \min_{\theta \in \mathbb{R}^p} L_T(\theta, \lambda^{(k)})$ .

    Construct matrix  $U^{(k)}$ , an orthonormal basis of  $\Omega^{L_T(\cdot, \lambda)}(\hat{\theta}(\lambda^{(k)}))$ .

    Define the locally equivalent joint optimization problem

$$\begin{aligned} & \min_{\lambda \in \Lambda} L(\mathbf{y}_V, f_{U^{(k)}\hat{\beta}(\lambda)}(\mathbf{X}_V)) \\ \text{s.t. } & \hat{\beta}(\lambda) = \arg \min_{\beta} L(\mathbf{y}_T, f_{U^{(k)}\beta}(\mathbf{X}_T)) + \sum_{i=1}^J \lambda_i P_i(U^{(k)}\beta) \end{aligned} \quad (35)$$

    Calculate  $\frac{\partial}{\partial \lambda} \hat{\beta}(\lambda)|_{\lambda=\lambda^{(k)}}$  where

$$\frac{\partial}{\partial \lambda} \hat{\beta}(\lambda) = - \left[ U^{(k)} \nabla^2 \left( L(\mathbf{y}_T, f_{U^{(k)}\beta}(\mathbf{X}_T)) + \sum_{i=1}^J \lambda_i P_i(U^{(k)}\beta) \right) \Big|_{\beta=\hat{\beta}(\lambda)} \right]^{-1} \left[ U^{(k)} \nabla P(U^{(k)}\beta)|_{\beta=\hat{\beta}(\lambda)} \right] \quad (36)$$

    with  $U^{(k)} \nabla^2$  and  $U^{(k)} \nabla$  are as defined in (12).

    Calculate the gradient  $\nabla_{\lambda} L(\mathbf{y}_V, f_{\hat{\theta}(\lambda)}(\mathbf{X}_V))|_{\lambda=\lambda^{(k)}}$  where

$$\nabla_{\lambda} L(\mathbf{y}_V, f_{\hat{\theta}(\lambda)}(\mathbf{X}_V)) = \left[ U^{(k)} \frac{\partial}{\partial \lambda} \hat{\beta}(\lambda) \right]^{\top} \left[ U^{(k)} \nabla L(\mathbf{y}_V, f_{U^{(k)}\beta}(\mathbf{X}_V)) \Big|_{\beta=\hat{\beta}(\lambda)} \right] \quad (37)$$

    Perform Neterov's update with step size  $t^{(k)}$ :

$$\begin{aligned} \eta &:= \lambda^{(k)} + \frac{k-1}{k+2} (\lambda^{(k)} - \lambda^{(k-1)}) \\ \lambda^{(k+1)} &:= \eta - t^{(k)} \nabla_{\lambda} L(\mathbf{y}_V, f_{\hat{\theta}(\lambda)}(\mathbf{X}_V)) \Big|_{\lambda=\eta} \end{aligned} \quad (38)$$

**if** the stopping criteria is reached or

$$L(\mathbf{y}_V, f_{\hat{\theta}(\lambda^{(k+1)})}(\mathbf{X}_V)) > L(\mathbf{y}_V, f_{\hat{\theta}(\lambda^{(k)})}(\mathbf{X}_V)), \quad (39)$$

**then**

    set  $\lambda^{(0)} := \lambda^{(k)}$  and break

**end if**

**end for**

**end while**

**return**  $\lambda^{(0)}$  and  $\hat{\theta}(\lambda^{(0)})$

---