

1 Appendix

1.1 Proof of Theorem 1

Proof. We will show that for a given λ_0 that satisfies the given conditions, the validation loss is continuously differentiable within some neighborhood of λ_0 . It then follows that if the theorem conditions hold true for almost every λ , then the validation loss is continuously differentiable with respect to λ at almost every λ .

Suppose the theorem conditions are satisfied at λ_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} :

$$\tilde{L}_T(\boldsymbol{\theta}, \lambda) =_{\mathbf{B}} \nabla L_T(\cdot, \lambda)|_{\boldsymbol{\theta}} \quad (1)$$

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 $\hat{\boldsymbol{\theta}}(\lambda_0)$:

$$_{\mathbf{B}} \nabla L_T(\cdot, \lambda_0)|_{\hat{\boldsymbol{\theta}}(\lambda_0)} = \tilde{L}_T(\hat{\boldsymbol{\theta}}(\lambda_0), \lambda_0) = 0 \quad (2)$$

From our assumptions, we know that there exists a neighborhood W containing λ_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 λ_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

$$\tilde{L}_T(\gamma(\lambda), \lambda) = \nabla_{\mathbf{B}} L_T(\cdot, \lambda)|_{\gamma(\lambda)} = 0 \quad (3)$$

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 $\lambda \in U$,

$$\hat{\boldsymbol{\theta}}(\lambda) = \arg \min_{\boldsymbol{\theta}} L_T(\boldsymbol{\theta}, \lambda) = \arg \min_{\boldsymbol{\theta} \in \Omega^{L_T}(\hat{\boldsymbol{\theta}}(\lambda_0), \lambda_0)} L_T(\boldsymbol{\theta}, \lambda) = \gamma(\lambda) \quad (4)$$

Therefore, we have shown that if λ_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 λ_0 . We can then apply the chain rule to get the gradient of the validation loss. \square

1.2 Gradient Derivations

1.2.1 Un-pooled Sparse Group Lasso

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

$$\begin{aligned} & \min_{\lambda \in \mathbb{R}_+^2} \frac{1}{2n} \left\| \mathbf{y}_V - \mathbf{X}_V \hat{\boldsymbol{\theta}}(\lambda) \right\|_2^2 \\ \text{s.t. } & \hat{\boldsymbol{\theta}}(\lambda) = \arg \min_{\boldsymbol{\theta}} \frac{1}{2n} \left\| \mathbf{y}_T - \mathbf{X}_T \boldsymbol{\theta} \right\|_2^2 + \sum_{m=1}^M \lambda_1^{(m)} \left\| \boldsymbol{\theta}^{(m)} \right\|_2 + \lambda_2 \left\| \boldsymbol{\theta} \right\|_1 + \frac{1}{2} \epsilon \left\| \boldsymbol{\theta} \right\|_2^2 \end{aligned} \quad (5)$$

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

$$\frac{1}{n} \mathbf{X}_{T, I(\boldsymbol{\lambda})}^\top \mathbf{X}_{T, I(\boldsymbol{\lambda})} + \mathbf{B}(\boldsymbol{\lambda}) + \epsilon \mathbf{I} \quad (6)$$

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

$$\frac{\lambda_1^{(m)}}{\|\boldsymbol{\theta}^{(m)}\|_2} \left(\mathbf{I} - \frac{1}{\|\boldsymbol{\theta}^{(m)}\|_2^2} \boldsymbol{\theta}^{(m)} \boldsymbol{\theta}^{(m)\top} \right) \quad (7)$$

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

$$\begin{aligned} & \min_{\boldsymbol{\lambda} \in \mathbb{R}_+^2} \frac{1}{2n} \left\| \mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \right\|_2^2 \\ \text{s.t. } & \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\beta}} \frac{1}{2n} \left\| \mathbf{y}_T - \mathbf{X}_{T, I(\boldsymbol{\lambda})} \boldsymbol{\beta} \right\|_2^2 + \sum_{m=1}^M \lambda_1^{(m)} \|\boldsymbol{\beta}^{(m)}\|_2 + \lambda_2 \|\boldsymbol{\beta}\|_1 + \frac{1}{2} \epsilon \|\boldsymbol{\beta}\|_2^2 \end{aligned} \quad (8)$$

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

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) &= \begin{bmatrix} \frac{\partial}{\partial \lambda_1^{(1)}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) & \dots & \frac{\partial}{\partial \lambda_1^{(M)}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) & \frac{\partial}{\partial \lambda_2} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \end{bmatrix} \\ &= - \left(\frac{1}{n} \mathbf{X}_{T, I(\boldsymbol{\lambda})}^\top \mathbf{X}_{T, I(\boldsymbol{\lambda})} + \mathbf{B}(\boldsymbol{\lambda}) + \epsilon \mathbf{I} \right)^{-1} \begin{bmatrix} \mathbf{C}(\hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) & \text{sgn}(\hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) \end{bmatrix} \end{aligned} \quad (9)$$

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

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ \frac{\hat{\boldsymbol{\beta}}^{(m)}(\boldsymbol{\lambda})}{\|\hat{\boldsymbol{\beta}}^{(m)}(\boldsymbol{\lambda})\|_2} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (10)$$

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

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{y}_V, \mathbf{X}_V \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) = \frac{1}{n} \left(\mathbf{X}_{V, I(\boldsymbol{\lambda})} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \right)^\top (\mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda})) \quad (11)$$

1.2.2 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 $\boldsymbol{\beta}$ and the H-P filter for the nonparametric estimates $\boldsymbol{\theta}$ is

$$\begin{aligned}
& \min_{\boldsymbol{\lambda} \in \mathbb{R}_+^2} \frac{1}{2} \left\| \mathbf{y}_V - \mathbf{X}_V \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) - (\mathbf{I} - \mathbf{I}_T) \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \right\|_2^2 \\
\text{s.t. } & \hat{\boldsymbol{\beta}}(\boldsymbol{\lambda}), \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\beta}, \boldsymbol{\theta}} \frac{1}{2} \left\| \mathbf{y}_T - \mathbf{X}_T \boldsymbol{\beta} - \mathbf{I}_T \boldsymbol{\theta} \right\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \frac{1}{2} \lambda_2 \|\boldsymbol{\beta}\|_2^2 + \frac{1}{2} \lambda_3 \|\mathbf{D}(\mathbf{z}) \boldsymbol{\theta}\|_2^2 + \frac{1}{2} \epsilon \|\boldsymbol{\theta}\|_2^2
\end{aligned} \tag{12}$$

The differentiable space is exactly the same as that given in Section 2.4.5. Also, all three conditions of Theorem 1 are satisfied. Note that the Hessian of the training criterion with respect to the basis in (54) is

$$H = \begin{bmatrix} \mathbf{I}_{I(\boldsymbol{\lambda})}^\top \mathbf{X}_T^\top \mathbf{X}_T \mathbf{I}_{I(\boldsymbol{\lambda})} + \lambda_2 \mathbf{I} & \mathbf{I}_{I(\boldsymbol{\lambda})}^\top \mathbf{X}_T^\top \mathbf{I}_T \\ \mathbf{I}_T^\top \mathbf{X}_T \mathbf{I}_{I(\boldsymbol{\lambda})} & \mathbf{I}_T^\top \mathbf{I}_T + \lambda_3 \mathbf{D}(\mathbf{z})^\top \mathbf{D}(\mathbf{z}) + \epsilon \mathbf{I} \end{bmatrix} \tag{13}$$

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

$$\begin{aligned}
& \min_{\boldsymbol{\lambda} \in \mathbb{R}_+^2} \frac{1}{2} \left\| \mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) - (\mathbf{I} - \mathbf{I}_T) \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \right\|_2^2 \\
\text{s.t. } & \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}), \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\eta}, \boldsymbol{\theta}} \frac{1}{2} \left\| \mathbf{y}_T - \mathbf{X}_{T, I(\boldsymbol{\lambda})} \boldsymbol{\eta} - \mathbf{I}_T \boldsymbol{\theta} \right\|_2^2 + \lambda_1 \|\boldsymbol{\eta}\|_1 + \frac{1}{2} \lambda_2 \|\boldsymbol{\eta}\|_2^2 + \frac{1}{2} \lambda_3 \|\mathbf{D}(\mathbf{z}) \boldsymbol{\theta}\|_2^2 + \frac{1}{2} \epsilon \|\boldsymbol{\theta}\|_2^2
\end{aligned} \tag{14}$$

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

$$\begin{bmatrix} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) \\ \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \lambda_1} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) & \frac{\partial}{\partial \lambda_3} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) & \frac{\partial}{\partial \lambda_3} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) \\ \frac{\partial}{\partial \lambda_1} \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) & \frac{\partial}{\partial \lambda_2} \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) & \frac{\partial}{\partial \lambda_3} \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \end{bmatrix} = -H^{-1} \begin{bmatrix} \text{sgn}(\hat{\boldsymbol{\eta}}(\boldsymbol{\lambda})) & \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{D}(\mathbf{z})^\top \mathbf{D}(\mathbf{z}) \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \end{bmatrix} \tag{15}$$

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

$$\nabla_{\boldsymbol{\lambda}} L_V(\boldsymbol{\lambda}) = - \left(\mathbf{X}_{V, I(\boldsymbol{\lambda})} \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) + (\mathbf{I} - \mathbf{I}_T) \frac{\partial}{\partial \boldsymbol{\lambda}} \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \right)^\top \left(\mathbf{y}_V - \mathbf{X}_{V, I(\boldsymbol{\lambda})} \hat{\boldsymbol{\eta}}(\boldsymbol{\lambda}) - (\mathbf{I} - \mathbf{I}_T) \hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}) \right) \tag{16}$$

1.3 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 Δx . 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 1 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

1.4 Joint Optimization with Accelerated Gradient Descent and Adaptive Restarts

Algorithm 2 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

$$\begin{aligned} & \min_{\lambda \in \Lambda} L(\mathbf{y}_V, f_{U^{(k)}\hat{\beta}(\lambda)}(\mathbf{X}_V)) \\ \text{s.t. } & \hat{\beta}(\lambda) = \arg \min_{\beta} L(\mathbf{y}_T, f_{U^{(k)}\beta}(\mathbf{X}_T)) + \sum_{i=1}^J \lambda_i P_i(U^{(k)}\beta) \end{aligned} \quad (17)$$

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

$$\frac{\partial}{\partial \lambda} \hat{\beta}(\lambda) = - \left[U^{(k)} \nabla^2 \left(L(\mathbf{y}_T, f_{U^{(k)}\beta}(\mathbf{X}_T)) + \sum_{i=1}^J \lambda_i P_i(U^{(k)}\beta) \right) \Big|_{\beta=\hat{\beta}(\lambda)} \right]^{-1} \left[U^{(k)} \nabla P(U^{(k)}\beta)|_{\beta=\hat{\beta}(\lambda)} \right] \quad (18)$$

 with $U^{(k)} \nabla^2$ and $U^{(k)} \nabla$ are as defined in (15).

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

$$\nabla_{\lambda} L(\mathbf{y}_V, f_{\hat{\theta}(\lambda)}(\mathbf{X}_V)) = \left[U^{(k)} \frac{\partial}{\partial \lambda} \hat{\beta}(\lambda) \right]^{\top} \left[U^{(k)} \nabla L(\mathbf{y}_V, f_{U^{(k)}\beta}(\mathbf{X}_V)) \Big|_{\beta=\hat{\beta}(\lambda)} \right] \quad (19)$$

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

$$\begin{aligned} \eta &:= \lambda^{(k)} + \frac{k-1}{k+2} (\lambda^{(k)} - \lambda^{(k-1)}) \\ \lambda^{(k+1)} &:= \eta - t^{(k)} \nabla_{\lambda} L(\mathbf{y}_V, f_{\hat{\theta}(\lambda)}(\mathbf{X}_V)) \Big|_{\lambda=\eta} \end{aligned} \quad (20)$$

if the stopping criteria is reached **or**

$$L(\mathbf{y}_V, f_{\hat{\theta}(\lambda^{(k+1)})}(\mathbf{X}_V)) > L(\mathbf{y}_V, f_{\hat{\theta}(\lambda^{(k)})}(\mathbf{X}_V)), \quad (21)$$

then

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

end if

end for

end while

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