

# MATH 680: Project Proposal

David Fleischer, Annik Gougeon

Last Update: 13 March, 2018

## Introduction

We wish to study the method of nonparametric regression with hierarchical penalization as outlined by Haris, Shojaie, and Simon (2016b). In particular, wish to investigate its properties, applications, as well as outline suitable classes of optimization techniques for such problems.

## Background

Consider the problem of estimating the relationship between responses  $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^n$  and predictors  $\mathbb{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times p}$ ,  $\mathbf{x}_i = [x_{i1}, \dots, x_{ip}] \in \mathbb{R}^p$ . Suppose that  $Y$  and  $\mathbb{X}$  are relation through an additive relationship

$$y_i = \sum_{j=1}^p f_j(x_{ij}) + \epsilon_i.$$

The method of estimation outlined herein focuses on a nonparametric method of basis expansion/projection estimators of  $Y$ . Specifically, for each predictor  $X_j = [x_{1j}, \dots, x_{nj}]^T \in \mathbb{R}^n$ , we generate the expansion of  $X_j$  according to a finite set of basis functions  $\{\psi_k(z)\}_{k=1}^K$ , where  $K$  is the truncation level that will be determined data-adaptively (discussed shortly). Let  $\Psi_K^{(j)} \in \mathbb{R}^{n \times K}$  be the set of basis functions correspond to predictor  $X_j$ , with entry  $(i, k)^{\text{th}}$  given by

$$\Psi_{K, (i, k)}^{(j)} = \psi_k(x_{ij}), \quad 1 \leq k \leq K, 1 \leq i \leq n.$$

Of present interest is the case of a *polynomial basis expansion*  $\psi_k^{(j)}(z) = z^k$  so that the  $j^{\text{th}}$  predictor undergoes the expansion

$$X_j = \begin{bmatrix} x_{1j} \\ \vdots \\ x_{nj} \end{bmatrix} \mapsto \Psi_K^{(j)} = \begin{bmatrix} \psi_1(x_{1j}) & \psi_2(x_{1j}) & \cdots & \psi_K(x_{1j}) \\ \vdots & \vdots & \ddots & \vdots \\ \psi_1(x_{nj}) & \psi_2(x_{nj}) & \cdots & \psi_K(x_{nj}) \end{bmatrix} = \begin{bmatrix} x_{1j} & x_{1j}^2 & \cdots & x_{1j}^K \\ \vdots & \vdots & \ddots & \vdots \\ x_{nj} & x_{nj}^2 & \cdots & x_{nj}^K \end{bmatrix}.$$

We may estimate the additive functions  $f_j(x_{ij})$  by the **sparse additive hierbasis** estimator  $\hat{f}_j(x_{ij}) = \sum_{k \leq K} \hat{\beta}_{j,k}^{\text{s-hier}} \psi_k(x_{ij})$ ,  $j = 1, \dots, p$ , such that each  $\hat{\beta}_j^{\text{s-hier}} \in \mathbb{R}^K$  is simultaneously estimated the penalized minimization problem

$$[\hat{\beta}_1^{\text{s-hier}}, \dots, \hat{\beta}_p^{\text{s-hier}}] = \arg \min_{\beta_j \in \mathbb{R}^K} \left\{ \frac{1}{2n} \left\| Y - \sum_{j=1}^p \Psi_K^{(j)} \beta_j \right\|_2^2 + \lambda \sum_{j=1}^p \Omega_j(\beta_j) + \frac{\lambda^2}{\sqrt{n}} \sum_{j=1}^p \left\| \Psi_K^{(j)} \beta_j \right\|_2 \right\}, \quad (1)$$

where

$$\Omega_j(\beta_j) = \frac{1}{\sqrt{n}} \sum_{k=1}^K w_k \left\| \Psi_{k:K}^{(j)} \beta_{j, k:K} \right\|_2$$

for  $w_k = k^m - (k-1)^m$ ,  $\Psi_{k:K}^{(j)}$  denotes the submatrix of columns  $k, k+1, \dots, K$ , and  $\beta_{k:K}$  the corresponding subvector of  $\beta$ .

The above penalty  $\Omega_j$  is designed to provide a data-driven method of truncating the basis complexity to some  $K_0 \leq K$ , derived from the hierarchical group lasso penalty (Zhao, Rocha, and Yu (2009)), leading to hierarchical sparsity of the fitted parameters  $\hat{\beta}_k = 0 \implies \hat{\beta}_{k'} = 0$ , for  $k' > k$ . The second penalization term  $\frac{\lambda^2}{n} \sum \left\| \Psi_K^{(j)} \beta_j \right\|_2$  imposes additional sparsity through a linked penalization constant  $\lambda^2$ .

## Solving for hierbasis Estimators

To solve (1) Haris, Shojaie, and Simon (2016b) applies the results of Zhao, Rocha, and Yu (2009), Jenatton et al. (2010), Jenatton et al. (2011). By writing the problem in the form

$$\min_{v \in \mathbb{R}^p} \left\{ \|u - v\|_2^2 + \lambda \Omega(v) \right\}, \quad (2)$$

where  $\Omega$  is a hierarchical penalty of the form described in Zhao, Rocha, and Yu (2009), we may apply an efficient proximal gradient descent algorithm with complexity  $O(p)$  (Jenatton et al. (2011)).

## Proposal

Of consideration for this project we wish to tackle the following questions:

- (1) Can the **hierbasis** estimator procedure offer a material gain over the lasso estimator (Tibshirani (1996))? Preliminary tests suggest a marginal sparsity improvement with no worse predictive power, but at the cost of computational complexity.
- (2) The **hierbasis** documentation (Haris, Shojaie, and Simon (2016a)) references a mixing parameter  $\alpha$  controlling the relative importance of the hierarchical and the sparsity-inducing penalties. How does the manipulation of this parameter affect its performance? Is it feasible to select  $\alpha$  through cross-validation?
- (3) What is the effect of changing the form of the weights  $w_k = k^m - (k-1)^m$  in the hierarchical penalty  $\Omega$ ?
- (4) What other optimization methods be used to solve **hierbasis**? Can other methods address the convergence issues?
- (5) Can the  $\ell_1$  norm be implemented in either/both penalties in order to induce stricter sparsity?

## References

- Haris, Asad, Ali Shojaie, and Noah Simon. 2016a. *HierBasis: Nonparametric Regression and Sparse Additive Modeling*.
- . 2016b. “Nonparametric Regression with Adaptive Truncation via a Convex Hierarchical Penalty.” *arXiv Preprint arXiv:1611.09972*.
- Jenatton, Rodolphe, Julien Mairal, Guillaume Obozinski, and Francis Bach. 2011. “Proximal Methods for Hierarchical Sparse Coding.” *Journal of Machine Learning Research* 12 (Jul): 2297–2334.
- Jenatton, Rodolphe, Julien Mairal, Guillaume Obozinski, and Francis R Bach. 2010. “Proximal Methods for Sparse Hierarchical Dictionary Learning.” In *ICML*, 487–94. 2010. Citeseer.
- Tibshirani, Robert. 1996. “Regression Shrinkage and Selection via the Lasso.” *Journal of the Royal Statistical Society. Series B (Methodological)*. JSTOR, 267–88.
- Zhao, Peng, Guilherme Rocha, and Bin Yu. 2009. “The Composite Absolute Penalties Family for Grouped and Hierarchical Variable Selection.” *The Annals of Statistics*. JSTOR, 3468–97.