

FAITHFUL VARIABLE SCREENING FOR HIGH-DIMENSIONAL CONVEX REGRESSION

BY MIN XU*, MINHUA CHEN[†] AND JOHN LAFFERTY[†]

We study the problem of variable selection in convex non-parametric regression. Under the assumption that the true regression function is convex and sparse, we develop a screening procedure to select a subset of variables that contains the relevant variables. Our approach is a two-stage quadratic programming method that estimates a sum of one-dimensional convex functions, followed by one-dimensional concave regression fits on the residuals. In contrast to previous methods for sparse additive models, the optimization is finite dimensional and requires no tuning parameters for smoothness. Under appropriate assumptions, we prove that the procedure is faithful in the population setting, yielding no false negatives. We give a finite sample statistical analysis, and introduce algorithms for efficiently carrying out the required quadratic programs. The approach leads to computational and statistical advantages over fitting a full model, and provides an effective, practical approach to variable screening in convex regression.

1. Introduction. Shape restrictions such as monotonicity, convexity, and concavity provide a natural way of limiting the complexity of many statistical estimation problems. Shape-constrained estimation is not as well understood as more traditional nonparametric estimation involving smoothness constraints. For instance, the minimax rate of convergence for multivariate convex regression has yet to be rigorously established in full generality. Even the one-dimensional case is challenging, and has been of recent interest ([Guntuboyina and Sen, 2013](#)).

In this paper we study the problem of variable selection in multivariate convex regression. Assuming that the regression function is convex and sparse, our goal is to identify the relevant variables. We show that it suffices to estimate a sum of one-dimensional convex functions, leading to significant computational and statistical advantages. This is in contrast to general nonparametric regression, where fitting an additive model can result in false negatives. Our approach is based on a two-stage quadratic programming procedure. In the first stage, we fit an convex additive

*Machine Learning Department, Carnegie Mellon University, Pittsburgh PA 15213

[†]Department of Statistics, University of Chicago, Chicago IL 60637

Keywords and phrases: nonparametric regression, convex regression, variable selection, quadratic programming, additive model

model, imposing a sparsity penalty. In the second stage, we fit a concave function on the residual for each variable. As we show, this non-intuitive second stage is in general necessary. Our first result is that this procedure is faithful in the population setting, meaning that it results in no false negatives, under mild assumptions on the density of the covariates. Our second result is a finite sample statistical analysis of the procedure, where we upper bound the statistical rate of variable screening consistency. An additional contribution is to show how the required quadratic programs can be formulated to be more scalable. We give simulations to illustrate our method, showing that it performs in a manner that is consistent with our analysis.

Estimation of convex functions arises naturally in several applications. Examples include geometric programming (Boyd and Vandenberghe, 2004), computed tomography (Prince and Willsky, 1990), target reconstruction (Lele, Kulkarni and Willsky, 1992), image analysis (Goldenshluger and Zeevi, 2006) and circuit design (Hannah and Dunson, 2012). Other applications include queuing theory (Chen and Yao, 2001) and economics, where it is of interest to estimate concave utility functions (Meyer and Pratt, 1968). See Lim and Glynn (2012) for other applications. Beyond cases where the assumption of convexity is natural, the convexity assumption can be attractive as a tractable, nonparametric relaxation of the linear model.

Recently, there has been increased research activity on shape-constrained estimation. Guntuboyina and Sen (2013) analyze univariate convex regression and show surprisingly that the risk of the MLE is adaptive to the complexity of the true function. Seijo and Sen (2011) and Lim and Glynn (2012) study maximum likelihood estimation of multivariate convex regression and independently establish its consistency. Cule, Samworth and Stewart (2010) and Kim and Samworth (2014) analyze log-concave density estimation and prove consistency of the MLE; the latter further show that log-concave density estimation has minimax risk lower bounded by $n^{-2/(d+1)}$ for $d \geq 2$, refuting a common notion that the condition of convexity is equivalent, in estimation difficulty, to the condition of having two bounded derivatives. Additive shape-constrained estimation has also been studied; Pya and Wood (2014) propose a penalized B-spline estimator while Chen and Samworth (2014) show the consistency of the MLE. To the best of our knowledge however, there has been no work on variable selection and estimation of high-dimensional convex functions.

Variable selection in general nonparametric regression or function estimation is a notoriously difficult problem. Lafferty and Wasserman (2008) develop a greedy procedure for adjusting bandwidths in a local linear regression estimator, and show that the procedure achieves the minimax rate as if the relevant variables were isolated in advance. But the method only provably scales to dimensions p that grow logarithmically in the sample size n , i.e., $p = O(\log n)$. This is in contrast to the high dimensional scaling behavior known to hold for sparsity selection in linear models using ℓ_1 penalization, where n is logarithmic in the dimension p . Bertin and Lecué (2008) develop an optimization-based approach in the nonparametric setting, apply-

ing the lasso in a local linear model at each test point. Here again, however, the method only scales as $p = O(\log n)$, the low-dimensional regime. An approximation theory approach to the same problem is presented in [DeVore, Petrova and Wojtaszczyk \(2011\)](#), using techniques based on hierarchical hashing schemes, similar to those used for “junta” problems ([Mossel, O’Donnell and Servedio, 2004](#)). Here it is shown that the sample complexity scales as $n > \log p$ if one adaptively selects the points on which the high-dimensional function is evaluated.

[Comminges and Dalalyan \(2012\)](#) show that the exponential scaling $n = O(\log p)$ is achievable if the underlying function is assumed to be smooth with respect to a Fourier basis. They also give support for the intrinsic difficulty of variable selection in nonparametric regression, giving lower bounds showing that consistent variable selection is not possible if $n < \log p$ or if $n < \exp s$, where s is the number of relevant variables. Variable selection over kernel classes is studied by [Koltchinskii and Yuan \(2010\)](#).

Perhaps more closely related to the present work is the framework studied by [Raskutti, Wainwright and Yu \(2012\)](#) for sparse additive models, where sparse regression is considered under an additive assumption, with each component function belonging to an RKHS. An advantage of working over an RKHS is that nonparametric regression with a sparsity-inducing regularization penalty can be formulated as a finite dimensional convex cone optimization. On the other hand, smoothing parameters for the component Hilbert spaces must be chosen, leading to extra tuning parameters that are difficult to select in practice. There has also been work on estimating sparse additive models over a spline basis, for instance the work of [Huang, Horowitz and Wei \(2010\)](#), but these approaches too require the tuning of smoothing parameters.

While nonparametric, the convex regression problem is naturally formulated using finite dimensional convex optimization, with no additional tuning parameters. The convex additive model can be used for convenience, without assuming it to actually hold, for the purpose of variable selection. As we show, our method scales to high dimensions, with a dependence on the intrinsic dimension s that scales polynomially, rather than exponentially as in the general case analyzed in [Comminges and Dalalyan \(2012\)](#).

In the following section we give a high-level summary of our technical results, including additive faithfulness, variable selection consistency, and high dimensional scaling. In [Section 3](#) we give a detailed account of our method and the conditions under which we can guarantee consistent variable selection. In [Section 4](#) we show how the required quadratic programs can be reformulated to be more efficient and scalable. In [Section 5](#) we give the details of our finite sample analysis, showing that a sample size growing as $n = O(\text{poly}(s) \log p)$ is sufficient for variable selection. In [Section 6](#) we report the results of simulations that illustrate our methods and theory. The full proofs are given in a technical appendix.

2. Overview of Results. In this section we provide a high-level description of our technical results. The full technical details, the precise statement of the results, and their detailed proofs are provided in following sections.

Our main contribution is an analysis of an additive approximation for identifying relevant variables in convex regression. We prove a result that shows when and how the additive approximation can be used without introducing false negatives in the population setting. In addition, we develop algorithms for the efficient implementation of the quadratic programs required by the procedure.

We first establish some notation, to be used throughout the paper. If \mathbf{x} is a vector, we use \mathbf{x}_{-k} to denote the vector with the k -th coordinate removed. If $\mathbf{v} \in \mathbb{R}^n$, then $v_{(1)}$ denotes the smallest coordinate of \mathbf{v} in magnitude, and $v_{(j)}$ denotes the j -th smallest; $\mathbf{1}_n \in \mathbb{R}^n$ is the all ones vector. If $X \in \mathbb{R}^p$ is a random variable and $S \subset \{1, \dots, p\}$, then X_S is the subvector of X restricted to the coordinates in S . Given n samples $X^{(1)}, \dots, X^{(n)}$, we use \bar{X} to denote the sample mean. Given a random variable X_k and a scalar x_k , we use $\mathbb{E}[\cdot | x_k]$ as a shorthand for $\mathbb{E}[\cdot | X_k = x_k]$.

2.1. Faithful screening. The starting point for our approach is the observation that least squares nonparametric estimation under convexity constraints is equivalent to a finite dimensional quadratic program. Specifically, the infinite dimensional optimization

$$(2.1) \quad \begin{aligned} & \text{minimize} && \sum_{i=1}^n (Y_i - f(\mathbf{x}_i))^2 \\ & \text{subject to} && f : \mathbb{R}^p \rightarrow \mathbb{R} \text{ is convex} \end{aligned}$$

is equivalent to the finite dimensional quadratic program

$$(2.2) \quad \begin{aligned} & \text{minimize}_{f, \beta} && \sum_{i=1}^n (Y_i - f_i)^2 \\ & \text{subject to} && f_j \geq f_i + \beta_i^T (\mathbf{x}_j - \mathbf{x}_i), \text{ for all } i, j. \end{aligned}$$

Here f_i is the estimated function value $f(\mathbf{x}_i)$, and the vectors $\beta_i \in \mathbb{R}^d$ represent supporting hyperplanes to the epigraph of f . See [Boyd and Vandenberghe \(2004\)](#), Section 6.5.5. Importantly, this finite dimensional quadratic program does not have tuning parameters for smoothing the function.

This formulation of convex regression is subject to the curse of dimensionality. Moreover, attempting to select variables by regularizing the subgradient vectors β_i with a group sparsity penalty is not effective. Intuitively, the reason is that all p components of the subgradient β_i appear in every convexity constraint $f_j \geq f_i + \beta_i^T (\mathbf{x}_j - \mathbf{x}_i)$; small changes to the subgradients may not violate the constraints. Experimentally, we find that regularization with a group sparsity penalty will make the subgradients of irrelevant variables small, but may not zero them out completely.

This motivates us to consider an additive approximation. Under a convex additive model, each component of the subgradient appears only in the convexity constraint for the corresponding variable:

$$(2.3) \quad f_{ki'} \geq f_{ki} + \beta_{ki}(x_{ki'} - x_{ki})$$

where $f_{ki} = f_k(x_{ki})$ and β_{ki} is the subgradient at point x_{ki} . As we show, this leads to an effective variable selection procedure. The shape constraints play an essential role. For general regression, using an additive approximation for variable selection may make errors. In particular, the nonlinearities in the regression function may result in an additive component being wrongly zeroed out. We show that this cannot happen for convex regression under appropriate conditions.

We say that a differentiable function f depends on variable x_k if $\partial_{x_k} f \neq 0$ with probability greater than zero. An additive approximation is given by

$$(2.4) \quad \{f_k^*\}, \mu^* := \arg \min_{f_1, \dots, f_p, \mu} \left\{ \mathbb{E} \left(f(X) - \mu - \sum_{k=1}^p f_k(X_k) \right)^2 : \mathbb{E} f_k(X_k) = 0 \right\}.$$

We say that f is *additively faithful* in case $f_k^* = 0$ implies that f does not depend on coordinate k . Additive faithfulness is a desirable property since it implies that an additive approximation may allow us to screen out irrelevant variables.

Our first result shows that convex multivariate functions are additively faithful under the following assumption on the distribution of the data.

DEFINITION 2.1. Let $p(\mathbf{x})$ be a density supported on $[0, 1]^p$. Then p satisfies the *boundary flatness condition* if for all j , and for all \mathbf{x}_{-j} ,

$$\frac{\partial p(\mathbf{x}_{-j} | x_j)}{\partial x_j} = \frac{\partial^2 p(\mathbf{x}_{-j} | x_j)}{\partial x_j^2} = 0 \quad \text{at } x_j = 0 \text{ and } x_j = 1.$$

As discussed in Section 3, this is a relatively weak condition. Our first result is that this condition suffices in the population setting of convex regression.

THEOREM 1. *Let p be a positive density supported on $C = [0, 1]^p$ that satisfies the boundary flatness property. If f is convex and twice differentiable, then f is additively faithful under p .*

Intuitively, an additive approximation zeroes out variable k when, fixing x_k , every “slice” of f integrates to zero. We prove this result by showing that “slices” of convex functions that integrate to zero cannot be “glued together” while still maintaining convexity.

While this shows that convex functions are additively faithful, it is difficult to estimate the optimal additive functions. The difficulty is that f_k^* need not be a convex function, as we show through a counterexample in Section 3. It may be possible to

estimate f_k^* with smoothing parameters, but, for the purpose of variable screening, it is sufficient in fact to approximate f_k^* by a *convex* additive model.

Our next result states that a convex additive fit, combined with a series of univariate concave fits, is faithful. We abuse notation in Theorem 2 and let the notation f_k^* represent convex additive components.

THEOREM 2. *Suppose $p(\mathbf{x})$ is a positive density on $C = [0, 1]^p$ that satisfies the boundary flatness condition. Suppose that f is convex and twice-differentiable. and that $\partial_{x_k} f$, $\partial_{x_k} p(\mathbf{x}_{-k} | x_k)$, and $\partial_{x_k}^2 p(\mathbf{x}_{-k} | x_k)$ are all continuous as functions on C . Define*

$$(2.5) \quad \{f_k^*\}_{k=1}^p, \mu^* = \arg \min_{\{f_k\}, \mu} \left\{ \mathbb{E} \left(f(X) - \mu - \sum_{k=1}^s f_k(X_k) \right)^2 : f_k \in \mathcal{C}^1, \mathbb{E} f_k(X_k) = 0 \right\}$$

where \mathcal{C}^1 is the set of univariate convex functions, and, with respect to f_k^* 's from above, define

$$(2.6) \quad g_k^* = \arg \min_{g_k} \left\{ \mathbb{E} \left(f(X) - \mu^* - \sum_{k' \neq k} f_{k'}^*(X_{k'}) - g_k \right)^2 : g_k \in -\mathcal{C}^1, \mathbb{E} g_k(X_k) = 0 \right\},$$

with $-\mathcal{C}^1$ denoting the set of univariate concave functions. Then $f_k^* = 0$ and $g_k^* = 0$ implies that f does not depend on x_k , i.e., $\partial_{x_k} f(\mathbf{x}) = 0$ with probability one.

This result naturally suggests a two-stage screening procedure for variable selection. In the first stage we fit a sparse convex additive model $\{\hat{f}_k\}$. In the second stage we fit a concave function \hat{g}_k to the residual for each variable having a zero convex component \hat{f}_k . If both $\hat{f}_k = 0$ and $\hat{g}_k = 0$, we can safely discard variable x_k . As a shorthand, we refer to this two-stage procedure as AC/DC. In the AC stage we fit an additive convex model. In the DC stage we fit decoupled concave functions on the residuals. The decoupled nature of the DC stage allows all of the fits to be carried out in parallel. The entire process involves no smoothing parameters. Our next result concerns the required optimizations, and their finite sample statistical performance.

2.2. Optimization. In Section 4 we present optimization algorithms for the additive convex regression stage. The convex constraints for the additive functions, analogous to the multivariate constraints (2.2), are that each component $f_k(\cdot)$ can be represented by its supporting hyperplanes, i.e.,

$$(2.7) \quad f_{ki'} \geq f_{ki} + \beta_{ki}(x_{ki'} - x_{ki}) \quad \text{for all } i, i'$$

where $f_{ki} := f_k(x_{ki})$ and β_{ki} is the subgradient at point x_{ki} . While this apparently requires $O(n^2 p)$ equations to impose the supporting hyperplane constraints, in fact, only $O(np)$ constraints suffice. This is because univariate convex functions are characterized by the condition that the subgradient, which is a scalar, must increase

monotonically. This observation leads to a reduced quadratic program with $O(np)$ variables and $O(np)$ constraints.

Directly applying a QP solver to this optimization is still computationally expensive for relatively large n and p . We thus develop a block coordinate descent method, where in each step we solve a sparse quadratic program involving $O(n)$ variables and $O(n)$ constraints. This is efficiently solved using optimization packages such as MOSEK. The details of these optimizations are given in Section 4.

2.3. Finite sample analysis. In Section 5 we analyze the finite sample variable selection consistency of convex additive modeling, without assuming that the true regression function f_0 is additive. Our analysis first establishes a sufficient deterministic condition for variable selection consistency, and then considers a stochastic setting. Our proof technique decomposes the KKT conditions for the optimization in a manner that is similar to the now standard *primal-dual witness* method (Wainwright, 2009).

We prove separate results that allow us to analyze false negative rates and false positive rates. To control false positives, we analyze scaling conditions on the regularization parameter λ_n for group sparsity needed to zero out irrelevant variables $k \in S^c$, where $S \subset \{1, \dots, p\}$ is the set of variables selected by the AC/DC algorithm in the population setting. To control false negatives, we analyze the restricted regression where the variables in S^c are zeroed out, following the primal-dual strategy.

Each of our theorems uses a subset of the following assumptions:

A1: X_S, X_{S^c} are independent.

A2: f_0 is convex and twice-differentiable. $\mathbb{E}f_0(X) = 0$.

A3: $\|f_0\|_\infty \leq sB$ and $\|f_k^*\| \leq B$ for all k .

A4: The noise is mean-zero sub-Gaussian, independent of X .

A5: The density $p(\mathbf{x})$ is bounded away from 0 and satisfies the boundary flatness condition.

In Assumption A3, $f^* = \sum_k f_k^*$ denotes the optimal additive projection of f_0 in the population setting.

Our analysis involves parameters α_+ and α_- , which are measures of the signal strength of the weakest variable:

$$\begin{aligned}\alpha_+ &= \inf_{f \in \mathcal{C}^p : \text{supp}(f) \subsetneq \text{supp}(f^*)} \left\{ \mathbb{E}(f_0(X) - f(X))^2 - \mathbb{E}(f_0(X) - f^*(X))^2 \right\} \\ \alpha_- &= \min_{k \in S : g_k^* \neq 0} \left\{ \mathbb{E}(f_0(X) - f^*(X))^2 - \mathbb{E}(f_0(X) - f^*(X) - g_k^*(X_k))^2 \right\}.\end{aligned}$$

Intuitively, if α_+ is small, then it is easier to make a false omission in the additive convex stage of the procedure. If α_- is small, then it is easier to make a false omission in the decoupled concave stage of the procedure.

We make strong assumptions on the covariates in A1 in order to make very weak assumptions on the true regression function f_0 in A2; in particular, we do not assume that f_0 is additive. Relaxing this condition is an important direction for future work. We also include an extra boundedness constraint to use new bracketing number results (Kim and Samworth, 2014).

Our main result is the following.

THEOREM 3. *Suppose assumptions A1-A5 hold. Let $\{\hat{f}_i\}$ be any AC solution and let $\{\hat{g}_k\}$ be any DC solution, both estimated with regularization parameter λ scaling as $\lambda = \Theta\left(sB\sqrt{\frac{1}{n}\log^2 np}\right)$. Suppose in addition that*

$$(2.8) \quad \alpha_f/\tilde{\sigma} \geq cB^2\sqrt{\frac{s^5}{n^{4/5}}\log^2 np}$$

$$(2.9) \quad \alpha_g^2/\tilde{\sigma} \geq cB^4\sqrt{\frac{s^5}{n^{4/5}}\log^2 2np}.$$

where $\tilde{\sigma} \equiv \max(\sigma, B)$ and c is a constant dependent only on b, c_1 .

Then, for sufficiently large n , with probability at least $1 - \frac{1}{n}$:

$$\begin{aligned} \hat{f}_k &\neq 0 \text{ or } \hat{g}_k \neq 0 \text{ for all } k \in S \\ \hat{f}_k &= 0 \text{ and } \hat{g}_k = 0 \text{ for all } k \notin S. \end{aligned}$$

This shows that variable selection consistency is achievable under exponential scaling of the ambient dimension, $p = O(\exp(n^c))$ for some $0 < c < 1$, as for linear models. The cost of nonparametric estimation is reflected in the scaling with respect to $s = |S|$, which can grow only as $o(n^{4/25})$.

We remark that Comminges and Dalalyan (2012) show that, even with the product distribution, under traditional smoothness constraints, variable selection is achievable only if $n > O(e^s)$. Here we demonstrate that convexity yields the scaling $n = O(\text{poly}(s))$.

3. Additive Faithfulness. For a general regression function, an additive approximation may result in a relevant variable being incorrectly marked as irrelevant. Such mistakes are inherent to the approximation and may persist even in the population setting. In this section we give examples of this phenomenon, and then show how the convexity assumption changes the behavior of the additive approximation. We work with $C = [0, 1]^p$ as the support of the distribution in this section but all of our results apply to general hypercubes. We begin with a lemma that characterizes the components of the additive approximation under mild conditions.

LEMMA 3.1. *Let F be a distribution on $C = [0, 1]^p$ with a positive density function p . Let $f : C \rightarrow \mathbb{R}$ be an integrable function, and define*

$$f_1^*, \dots, f_p^*, \mu^* := \arg \min \left\{ \mathbb{E} \left(f(X) - \mu - \sum_{k=1}^p f_k(X_k) \right)^2 : \mathbb{E} f_k(X_k) = 0, \forall k = 1, \dots, p \right\}.$$

Then $\mu^* = \mathbb{E} f(X)$,

$$(3.1) \quad f_k^*(x_k) = \mathbb{E} \left[f(X) - \sum_{k' \neq k} f_{k'}^*(X_{k'}) \mid x_k \right] - \mathbb{E} f(X),$$

and this solution is unique.

Lemma 3.1 follows from the stationarity conditions of the optimal solution. This result is known, and criterion (3.1) is used in the backfitting algorithm for fitting additive models. We include a proof as our results build on it.

PROOF. Let $f_1^*, \dots, f_p^*, \mu^*$ be the minimizers as defined. We first show that the optimal μ is $\mu^* = \mathbb{E} f(X)$ for any f_1, \dots, f_k such that $\mathbb{E} f_k(X_k) = 0$. This follows from the stationarity condition, which states that $\mu^* = \mathbb{E}[f(X) - \sum_k f_k(X_k)] = \mathbb{E}[f(X)]$. Uniqueness is apparent because the second derivative is strictly larger than zero and strong convexity is guaranteed.

We now turn our attention toward the f_k^* s. It must be that f_k^* minimizes

$$(3.2) \quad \mathbb{E} \left[\left(f(X) - \mu^* - \sum_{k' \neq k} f_{k'}^*(X_{k'}) - f_k(X_k) \right)^2 \right]$$

subject to $\mathbb{E} f_k(X_k) = 0$. Fixing x_k , we will show that the value

$$(3.3) \quad \mathbb{E}[f(X) - \sum_{k' \neq k} f_{k'}^*(X_{k'}) \mid x_k] - \mu^*$$

uniquely minimizes

$$(3.4) \quad \min_{f_k(x_k)} \int_{\mathbf{x}_{-k}} p(\mathbf{x}) \left(f(\mathbf{x}) - \sum_{k' \neq k} f_{k'}^*(x_{k'}) - f_k(x_k) - \mu^* \right)^2 d\mathbf{x}_{-k}.$$

The first-order optimality condition gives us

$$(3.5) \quad \int_{\mathbf{x}_{-k}} p(\mathbf{x}) f_k(x_k) d\mathbf{x}_{-k} = \int_{\mathbf{x}_{-k}} p(\mathbf{x}) \left(f(\mathbf{x}) - \sum_{k' \neq k} f_{k'}^*(x_{k'}) - \mu^* \right) d\mathbf{x}_{-k}$$

$$(3.6) \quad p(x_k) f_k(x_k) = \int_{\mathbf{x}_{-k}} p(x_k) p(\mathbf{x}_{-k} \mid x_k) \left(f(\mathbf{x}) - \sum_{k' \neq k} f_{k'}^*(x_{k'}) - \mu^* \right) d\mathbf{x}_{-k}$$

$$(3.7) \quad f_k(x_k) = \int_{\mathbf{x}_{-k}} p(\mathbf{x}_{-k} \mid x_k) \left(f(\mathbf{x}) - \sum_{k' \neq k} f_{k'}^*(x_{k'}) - \mu^* \right) d\mathbf{x}_{-k}$$

The square error objective is strongly convex, and the second derivative with respect to $f_k(x_k)$ is $2p(x_k)$, which is always positive under the assumption that p is positive. Therefore, the solution $f_k^*(x_k) = \mathbb{E}[f(X) | x_k] - \mathbb{E}f(X)$ is unique. Noting that $\mathbb{E}[f(X) - \sum_{k' \neq k} f_{k'}^*(X_{k'}) | x_k] - \mathbb{E}f(X)$ has mean zero as a function of x_k completes the proof. \square

In the case that the distribution in Lemma 3.1 is a product distribution, the additive components take on a simple form.

COROLLARY 3.1. *Let F be a product distribution on $C = [0, 1]^p$ with density function p which is positive on C . Let $\mu^*, f_k^*(x_k)$ be defined as in Lemma 3.1. Then $\mu^* = \mathbb{E}f(X)$ and $f_k^*(x_k) = \mathbb{E}[f(X) | x_k] - \mathbb{E}f(X)$ and this solution is unique.*

In particular, if F is the uniform distribution, then $f_k^*(x_k) = \int f(x_k, \mathbf{x}_{-k}) d\mathbf{x}_{-k}$.

EXAMPLE 3.1. Using Corollary 3.1, we give two examples of *additively unfaithfulness* under the uniform distribution—where relevant variables are erroneously marked as irrelevant under an additive approximation. First, consider the following function:

$$(3.8) \quad f(x_1, x_2) = \sin(2\pi x_1) \sin(2\pi x_2) \quad (\text{egg carton})$$

defined for $(x_1, x_2) \in [0, 1]^2$. Then $\int_{x_2} f(x_1, x_2) dx_2 = 0$ and $\int_{x_1} f(x_1, x_2) dx_1 = 0$ for each x_1 and x_2 . An additive approximation would set $f_1 = 0$ and $f_2 = 0$. Next, consider the function

$$(3.9) \quad f(x_1, x_2) = x_1 x_2 \quad (\text{tilting slope})$$

defined for $x_1 \in [-1, 1]$, $x_2 \in [0, 1]$. In this case $\int_{x_1} f(x_1, x_2) dx_1 = 0$ for each x_2 ; therefore, we expect $f_2 = 0$ under the additive approximation. This function, for every fixed x_2 , is a zero-intercept linear function of x_1 with slope x_2 .

In order to exploit additive models in variable selection, it is important to understand when the additive approximation accurately captures all of the relevant variables. We call this property *additive faithfulness*. We first formalize the concept that a multivariate function f *does not depend on* a coordinate x_k .

DEFINITION 3.1. Let $C = [0, 1]^p$ and let $f : C \rightarrow \mathbb{R}$. We say that f *does not depend on coordinate k* if for all \mathbf{x}_{-k} , $f(x_k, \mathbf{x}_{-k})$ is a constant as a function of x_k . If f is differentiable, then f does not depend on k if $\partial_{x_k} f(x_k, \mathbf{x}_{-k})$ is 0 for all \mathbf{x}_{-k} .

In addition, suppose we have a distribution over C and the additive approximation

$$(3.10) \quad f_k^*, \mu^* := \arg \min_{f_1, \dots, f_p, \mu} \left\{ \mathbb{E} \left[\left(f(X) - \sum_{k=1}^p f_k(X_k) - \mu \right)^2 \right] : \mathbb{E} f_k(X_k) = 0 \right\}.$$

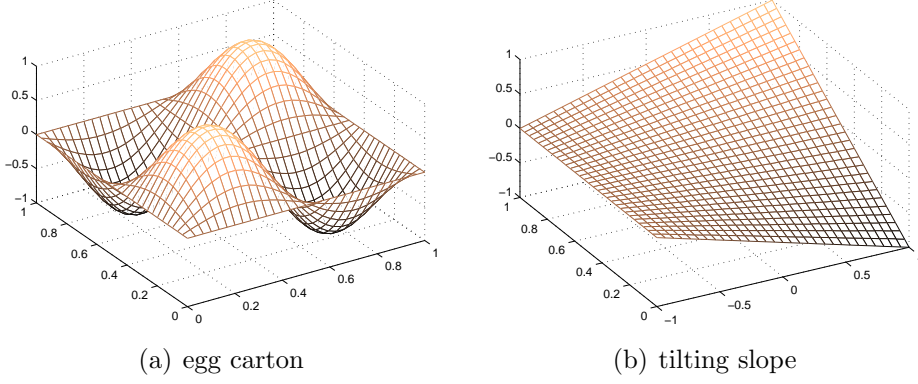


FIG 1. *Two additively unfaithful functions. Relevant variables are zeroed out under an additive approximation because every “slice” of the function integrates to zero.*

We say that f is *additively faithful* under F if $f_k^* = 0$ implies that f does not depend on coordinate k .

Additive faithfulness is an attractive property because it implies that, in the population setting, the additive approximation yields a consistent variable screening.

3.1. Additive Faithfulness of Convex Functions. We now show that under a general class of distributions which we characterize below, convex multivariate functions are additively faithful.

DEFINITION 3.2. Let $p(\mathbf{x})$ be a density supported on $[0, 1]^p$. $p(\mathbf{x})$ satisfies the *boundary flatness condition* if for all j , and for all \mathbf{x}_{-j} ,

$$(3.11) \quad \frac{\partial p(\mathbf{x}_{-j} | x_j)}{\partial x_j} = \frac{\partial^2 p(\mathbf{x}_{-j} | x_j)}{\partial x_j^2} = 0 \quad \text{at } x_j = 0, x_j = 1$$

The boundary flatness condition is a weak condition. For instance, it is satisfied when the density is flat at the boundary of support—more precisely, when the joint density satisfies the property that $\frac{\partial p(x_j, \mathbf{x}_{-j})}{\partial x_j} = \frac{\partial^2 p(x_j, \mathbf{x}_{-j})}{\partial x_j^2} = 0$ at boundary points $x_j = 0, x_j = 1$. The boundary flatness property is also trivially satisfied when p is a product density.

The following theorem is the main result of this section.

THEOREM 3.1. *Let p be a positive density supported on $C = [0, 1]^p$ that satisfies the boundary flatness property. If f is convex and twice differentiable, then f is additively faithful under p .*

We pause to give some intuition before we presenting the full proof. Suppose that the underlying distribution has a product density. Then we know from Lemma 3.1 that the additive approximation zeroes out k when, fixing x_k , every “slice” of f

integrates to zero. We prove Theorem 3.1 by showing that “slices” of convex functions that integrate to zero cannot be “glued together” while still maintaining convexity.

PROOF. Fixing k and using the result of Lemma 3.1, we need only show that for all x_k , $\mathbb{E}[f(X) - \sum_{k'} f_{k'}(X_{k'}) | x_k] - \mathbb{E}f(X) = 0$ implies that f does not depend on coordinate k , i.e., $\partial_{x_k} f(\mathbf{x}) = 0$ for all \mathbf{x} .

Let us use the shorthand notation that $r(\mathbf{x}_{-k}) = \sum_{k' \neq k} f_{k'}(x_{k'})$ and assume without loss of generality that $\mu^* = E[f(X)] = 0$. We then assume that for all x_k ,

$$(3.12) \quad \mathbb{E}[f(X) - r(X_{-k}) | x_k] \equiv \int_{\mathbf{x}_{-k}} p(\mathbf{x}_{-k} | x_k) (f(\mathbf{x}) - r(\mathbf{x}_{-k})) d\mathbf{x}_{-k} = 0.$$

We let $p'(\mathbf{x}_{-k} | x_k)$ denote $\frac{\partial p(\mathbf{x}_{-k} | x_k)}{\partial x_k}$ and $p''(\mathbf{x}_{-k} | x_k)$ denote $\frac{\partial^2 p(\mathbf{x}_{-k} | x_k)}{\partial x_k^2}$ and likewise for $f'(x_k, \mathbf{x}_{-k})$ and $f''(x_k, \mathbf{x}_{-k})$. We then differentiate under the integral, valid because all functions are bounded, and obtain

$$(3.13) \quad \int_{\mathbf{x}_{-k}} p'(\mathbf{x}_{-k} | x_k) (f(\mathbf{x}) - r(\mathbf{x}_{-k})) + p(\mathbf{x}_{-k} | x_k) f'(x_k, \mathbf{x}_{-k}) d\mathbf{x}_{-k} = 0$$

$$(3.14) \quad \int_{\mathbf{x}_{-k}} p''(\mathbf{x}_{-k} | x_k) (f(\mathbf{x}) - r(\mathbf{x}_{-k})) + 2p'(\mathbf{x}_{-k} | x_k) f'(x_k, \mathbf{x}_{-k}) + p(\mathbf{x}_{-k} | x_k) f''(x_k, \mathbf{x}_{-k}) d\mathbf{x}_{-k} = 0.$$

By the boundary flatness condition, we have that $p''(\mathbf{x}_{-k} | x_k)$ and $p'(\mathbf{x}_{-k} | x_k)$ are zero at $x_k = x_k^0 \equiv 0$. The integral equations then reduce to the following:

$$(3.15) \quad \int_{\mathbf{x}_{-k}} p(\mathbf{x}_{-k} | x_k^0) f'(x_k^0, \mathbf{x}_{-k}) d\mathbf{x}_{-k} = 0$$

$$(3.16) \quad \int_{\mathbf{x}_{-k}} p(\mathbf{x}_{-k} | x_k^0) f''(x_k^0, \mathbf{x}_{-k}) d\mathbf{x}_{-k} = 0.$$

Because f is convex, $f(x_k, \mathbf{x}_{-k})$ must be a convex function of x_k for all \mathbf{x}_{-k} . Therefore, for all \mathbf{x}_{-k} , $f''(x_k^0, \mathbf{x}_{-k}) \geq 0$. Since $p(\mathbf{x}_{-k} | x_k^0) > 0$ by the assumption that p is a positive density, we have that $\forall \mathbf{x}_{-k}$, $f''(x_k^0, \mathbf{x}_{-k}) = 0$ necessarily.

The Hessian of f at (x_k^0, \mathbf{x}_{-k}) then has a zero at the k -th main diagonal entry. A positive semidefinite matrix with a zero on the k -th main diagonal entry must have only zeros on the k -th row and column; see proposition 7.1.10 of [Horn and Johnson \(1990\)](#). Thus, at all \mathbf{x}_{-k} , the gradient of $f'(x_k^0, \mathbf{x}_{-k})$ with respect to \mathbf{x}_{-k} must be zero. Therefore, $f'(x_k^0, \mathbf{x}_{-k})$ must be constant for all \mathbf{x}_{-k} . By equation 3.15, we conclude that $f'(x_k^0, \mathbf{x}_{-k}) = 0$ for all \mathbf{x}_{-k} . We can use the same reasoning for the case where $x_k = x_k^1$ and deduce that $f'(x_k^1, \mathbf{x}_{-k}) = 0$ for all \mathbf{x}_{-k} .

Because $f(x_k, \mathbf{x}_{-k})$ as a function of x_k is convex, it must be that, for all $x_k \in (0, 1)$ and for all \mathbf{x}_{-k} ,

$$(3.17) \quad 0 = f'(x_k^0, \mathbf{x}_{-k}) \leq f'(x_k, \mathbf{x}_{-k}) \leq f'(x_k^1, \mathbf{x}_{-k}) = 0$$

Therefore f does not depend on x_k . □

Theorem 3.1 plays an important role in our finite sample analysis, where we show that the additive approximation is variable screening consistent, even when the true function is not additive.

REMARK 3.1. We assume twice differentiability in Theorems 3.1 to simplify the proof. We expect, however, that this smoothness condition is not necessary—every convex function can be approximated arbitrarily well by a smooth convex function.

REMARK 3.2. We have not found natural conditions under which the opposite direction of additive faithfulness holds—conditions implying that if f does not depend on coordinate k , then f_k^* will be zero in the additive approximation. Suppose, for example, that f is only a function of X_1, X_2 , and that (X_1, X_2, X_3) follows a degenerate 3-dimensional distribution where $X_3 = f(X_1, X_2) - f^*(X_1) - f_2^*(X_2)$. In this case X_3 exactly captures the additive approximation error. The best additive approximation of f would have a component $f_3^*(x_3) = x_3$ even though f does not depend on x_3 .

REMARK 3.3. In Theorem 3.1, we do not assume a parametric form for the additive components. The additive approximations may not be faithful if we assume a parametric form for each of the components. For example, suppose we approximate a convex function $f(X)$ by a linear form $X\beta$. The optimal linear function in the population setting is $\beta^* = \Sigma^{-1}\text{Cov}(X, f(X))$. Suppose the X 's are independent and follow a symmetric distribution and suppose $f(\mathbf{x}) = x_1^2 - \mathbb{E}[X_1^2]$, then $\beta_1^* = \mathbb{E}[X_1 f(X)] = \mathbb{E}[X_1^3 - X_1 \mathbb{E}[X_1^2]] = 0$.

REMARK 3.4. It is possible to get a similar result for distributions with unbounded support, by using a limit condition $\lim_{|x_k| \rightarrow \infty} \frac{\partial p(\mathbf{x}_{-k} | x_k)}{\partial x_k} = 0$. Such a limit condition however is not obeyed by many common distributions such as the multivariate Gaussian distribution. The next example shows that certain convex functions are not additive faithful under certain multivariate Gaussian distributions.

EXAMPLE 3.2. Consider a two dimensional quadratic function $f(\mathbf{x}) = \mathbf{x}^\top H \mathbf{x} + c$ where $H = \begin{pmatrix} H_{11} & H_{12} \\ H_{12} & H_{22} \end{pmatrix}$ is positive definite and a Gaussian distribution $X \sim N(0, \Sigma)$ where $\Sigma = \begin{pmatrix} 1 & \alpha \\ \alpha & 1 \end{pmatrix}$. As we show in Section 9 of the Appendix, the additive approximation has the following closed form.

$$\begin{aligned} f_1^*(x_1) &= \left(\frac{T_1 - T_2 \alpha^2}{1 - \alpha^4} \right) x_1^2 + c_1 \\ f_2^*(x_2) &= \left(\frac{T_2 - T_1 \alpha^2}{1 - \alpha^4} \right) x_2^2 + c_2 \end{aligned}$$

Where $T_1 = H_{11} + 2H_{12}\alpha + H_{22}\alpha^2$, $T_2 = H_{22} + 2H_{12}\alpha + H_{11}\alpha^2$, c_1, c_2 are constants such that f_1^* and f_2^* both have mean zero. Let $H = \begin{pmatrix} 1.6 & 2 \\ 2 & 5 \end{pmatrix}$, then it is easy to check that if $\alpha = -\frac{1}{2}$, then $f_1^* = 0$ and additive faithfulness is violated, if $\alpha > \frac{1}{2}$, then f_1^* is a concave function. We take the setting where $\alpha = -0.5$, compute the optimal additive functions via numerical simulation, and show the results in Figure 2(a)– f_1^* is zero as expected.

Although the Gaussian distribution does not satisfy the boundary flatness condition, it is possible to approximate the Gaussian distribution arbitrarily well with distributions that do satisfy the boundary flatness conditions.

EXAMPLE 3.3. Let Σ be as in Example 3.2 with $\alpha = -0.5$ so that $f_1^* = 0$. Consider a mixture $\lambda U[-(b + \epsilon), b + \epsilon]^2 + (1 - \lambda)N_b(0, \Sigma)$ where $N_b(0, \Sigma)$ is the density of a *truncated* bivariate Gaussian bounded in $[-b, b]^2$ and $U[-(b + \epsilon), b + \epsilon]^2$ is the uniform distribution over a square. The uniform distribution is supported over a slightly larger square to satisfy the boundary flatness conditions.

When b is large, ϵ is small, and λ is small, the mixture closely approximates the Gaussian distribution but is still additively faithful for convex functions. Figure 2(b) shows the optimal additive components under the mixture distribution, computed by numerical integration with $b = 5, \epsilon = 0.3, \lambda = 0.0001$. True to our theory, f_1^* , which is zero under the Gaussian distribution, is nonzero under the mixture approximation to the Gaussian distribution. We note that the magnitude $\mathbb{E}f_1^*(X_1)^2$, although nonzero, is very small, consistent with the fact that the mixture distribution closely approximates the Gaussian distribution.

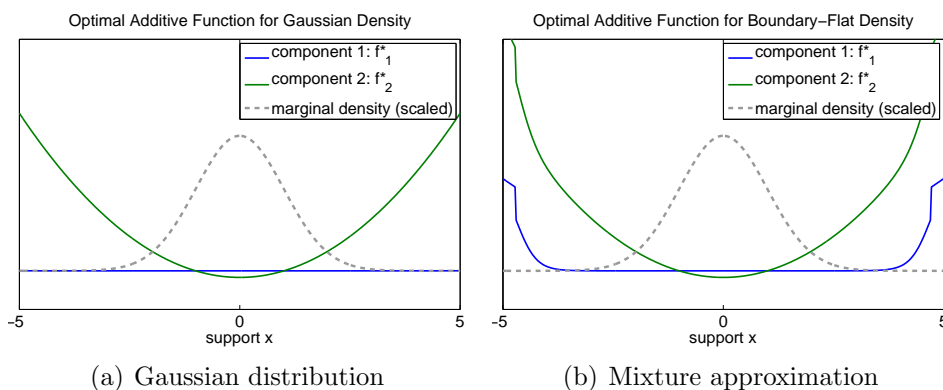


FIG 2. Optimal additive projection of the quadratic function described in Example 3.2 under both the Gaussian distribution described in Example 3.2 and under the approximately Gaussian mixture distribution described in Example 3.3. For the mixture approximation, we used $b = 5, \epsilon = 0.3, \lambda = 0.0001$ where the parameters are defined in Example 3.3. This example shows the effect and the importance of the boundary flatness conditions.

3.2. Convex Additive Models. Although convex functions are additively faithful—under appropriate conditions—it is difficult to estimate the optimal additive functions f_k^* s as defined in equation (3.10). The reason is that f_k^* need not be a convex function, as example 3.2 and example 3.3 show. It may be possible to estimate f_k^* via smoothing, but we prefer an approach that is free of smoothing parameters. Since the true regression function f is convex, we approximate the additive model with a *convex* additive model. We abuse notation and, for the rest of the paper, use the notation f_k^* to represent convex additive fits:

$$(3.18) \quad \{f_k^*\}_{k=1}^p = \arg \min \left\{ \mathbb{E} \left(f(X) - \sum_{k=1}^p f_k(X_k) \right)^2 : f_k \in \mathcal{C}^1, \mathbb{E} f_k(X_k) = 0 \right\}$$

where \mathcal{C}^1 is the set of univariate convex functions. The convex functions $\{f_k^*\}$ are not additively faithful by themselves, i.e., it could be that the true function f depends on variable k but $f_k^* = 0$. However, faithfulness can be restored by coupling the f_k^* 's with a set of univariate concave fits on the *residual* $f - f^*$:

$$(3.19) \quad g_k^* = \arg \min \left\{ \mathbb{E} \left(f(X) - \sum_{k' \neq k} f_{k'}^*(X_{k'}) - g_k \right)^2 : g_k \in -\mathcal{C}^1, \mathbb{E} g_k(X_k) = 0 \right\}.$$

THEOREM 3.2. *Suppose $p(\mathbf{x})$ is a positive density on $C = [0, 1]^p$ that satisfies the boundary flatness condition. Suppose that f is convex and twice differentiable, and that $\partial_{x_k} p(\mathbf{x}_{-k} | x_k)$, and $\partial_{x_k}^2 p(\mathbf{x}_{-k} | x_k)$ are all continuous as functions of x_k . Let f_k^* and g_k^* be as defined in equations (3.18) and (3.19), then the f_k^* 's and the g_k^* 's are unique. Furthermore, $f_k^* = 0$ and $g_k^* = 0$ implies that $\partial_{x_k} f(\mathbf{x}) = 0$, that is, f does not depend on x_k .*

Before we can prove the theorem, we need a lemma that generalizes Theorem 3.1.

LEMMA 3.2. *Suppose $p(x)$ is a positive density on $C = [0, 1]^p$ satisfying the boundary flatness condition. Let $f(\mathbf{x})$ be a convex twice differentiable function on C . Let $\phi(\mathbf{x}_{-k})$ be any function that does not depend on x_k . Then, we have that the unconstrained univariate function*

$$(3.20) \quad h_k^* = \arg \min_{f_k} \mathbb{E} \left[\left(f(X) - \phi(X_{-k}) - h_k(X_k) \right)^2 \right]$$

is given by $h_k^(x_k) = \mathbb{E}[f(X) - \phi(X_{-k}) | x_k]$, and $h_k^* = 0$ implies that $\partial_{x_k} f(\mathbf{x}) = 0$.*

PROOF. In the proof of Theorem 3.1, the only property of $r(\mathbf{x}_{-k})$ we used was the fact that $\partial_{x_k} r(\mathbf{x}_{-k}) = 0$. Therefore, the proof here is identical to that of Theorem 3.1 except that we replace $r(\mathbf{x}_{-k})$ with $\phi(\mathbf{x}_{-k})$. \square

PROOF OF THEOREM 3.2. Fix k . Let f_k^* and g_k^* be defined as in equation 3.18 and equation 3.19. Let $\phi(\mathbf{x}_{-k}) \equiv \sum_{k' \neq k} f_{k'}^*(x_{k'})$.

Then we have that

$$(3.21) \quad f_k^* = \arg \min_{f_k} \left\{ \mathbb{E} \left(f(X) - \sum_{k' \neq k} f_{k'}^*(X_{k'}) - f_k \right)^2 : f_k \in \mathcal{C}^1, \mathbb{E} f_k(X_k) = 0 \right\}$$

$$(3.22) \quad g_k^* = \arg \min_{g_k} \left\{ \mathbb{E} \left(f(X) - \sum_{k' \neq k} f_{k'}^*(X_{k'}) - g_k \right)^2 : g_k \in -\mathcal{C}^1, \mathbb{E} g_k(X_k) = 0 \right\}$$

Let us suppose that $f_k^* = g_k^* = 0$. It must be then that

$$\arg \min_{c \in \mathbb{R}} \mathbb{E} \left(f(X) - \phi(X_{-k}) - c(X_k^2 - m_k^2) \right)^2 = 0$$

where $m_k^2 \equiv \mathbb{E} X_k^2$; this is because $c(x_k^2 - m_k^2)$ is either convex or concave in x_k and it is centered, i.e. $\mathbb{E}[X_k^2 - m_k^2] = 0$. Since the optimum has a closed form $c^* = \frac{\mathbb{E}[(f(X) - \phi(X_{-k}))(X_k^2 - m_k^2)]}{\mathbb{E} X_k^2}$, we deduce that

$$\begin{aligned} \mathbb{E}[(f(X) - \phi(X_{-k}))(X_k^2 - m_k^2)] \\ = \mathbb{E}[(f(X) - \phi(X_{-k}))X_k^2] = \mathbb{E}[\mathbb{E}[f(X) - \phi(X_{-k}) | X_k]X_k^2] = 0 \end{aligned}$$

We denote $h_k^*(x_k) = \mathbb{E}[f(X) - \phi(X_{-k}) | x_k]$.

Under the derivative continuity conditions in the theorem, we apply Lemma 8.3 in the appendix and know that $h_k^*(x_k)$ is twice-differentiable and has a second derivative bounded away from $-\infty$. Therefore, for a large enough positive scalar α , $h_k^*(x_k) + \alpha(x_k^2 - m_k^2)$ has a non-negative second derivative and is thus convex.

Because we assumed $f^* = g^* = 0$, it must be that

$$\arg \min_{c \in \mathbb{R}} \mathbb{E} \left(f(X) - \phi(X_{-k}) - c(h_k^*(X_k) + \alpha(X_k^2 - m_k^2)) \right)^2 = 0$$

This is because $c(h_k^*(x_k) + \alpha(x_k^2 - m_k^2))$ is convex for $c \geq 0$ and concave for $c \leq 0$ and it is a centered function.

$$\text{Again, } c^* = \frac{\mathbb{E}[(f(X) - \phi(X_{-k}))(h_k^*(X_k) + \alpha(X_k^2 - m_k^2))]}{\mathbb{E}(h_k^*(X_k) + \alpha(X_k^2 - m_k^2))^2} = 0, \text{ so}$$

$$\begin{aligned} \mathbb{E}[(f(X) - \phi(X_{-k}))(h_k^*(X_k) + \alpha(X_k^2 - m_k^2))] &= \mathbb{E}[(f(X) - \phi(X_{-k}))h_k^*(X_k)] \\ &= \mathbb{E}[\mathbb{E}[f(X) - \phi(X_{-k}) | X_k]h_k^*(X_k)] \\ &= \mathbb{E}h_k^*(X_k)^2 = 0 \end{aligned}$$

where the first equality follows because $\mathbb{E}[(f(X) - \phi(X_{-k}))(X_k^2 - m_k^2)] = 0$. Therefore, we get that $h_k^* = 0$. Now we use Lemma 3.2 with $\phi(\mathbf{x}_{-k}) = f(\mathbf{x}) - \sum_{k' \neq k} f_{k'}^*(x_{k'})$ and conclude that $f_k^* = 0$ and $g_k^* = 0$ together imply that f does not depend on x_k .

Now we turn to uniqueness. Suppose for sake of contradiction that f^* and f' are optimal solutions to (3.18) and $\mathbb{E}(f' - f^*)^2 > 0$. $f^* + \lambda(f' - f^*)$ for any $\lambda \in [0, 1]$

must then also be an optimal solution by convexity of the objective and constraint. However, the second derivative of the objective $\mathbb{E}(f - f^* - \lambda(f' - f^*))^2$ with respect to λ is $2\mathbb{E}(f' - f^*)^2 > 0$. The objective is thus strongly convex and $\mathbb{E}(f^* - f')^2 = 0$. We now apply Lemma 8.4 by letting $\phi_k = f_k^* - f'_k$. We conclude that $\mathbb{E}(f_k^* - f'_k)^2 = 0$ for all k . The uniqueness of g^* is proved similarly. \square

3.3. Estimation Procedure. Theorem 3.2 naturally suggests a two-stage screening procedure for variable selection in the population setting. In the first stage, we fit a convex additive model.

$$(3.23) \quad f_1^*, \dots, f_p^* = \arg \min_{f_1, \dots, f_p \in \mathcal{C}_0^1} \mathbb{E} \left(f(X) - \sum_{k=1}^p f_k(X_k) \right)^2$$

where we denote \mathcal{C}_0^1 ($-\mathcal{C}_0^1$) as the set of one-dimensional convex (resp. concave) functions with population mean zero. In the second stage, for every variable marked as irrelevant in the first stage, we fit a univariate *concave* function separately on the residual for that variable. for each k such that $f_k^* = 0$:

$$(3.24) \quad g_k^* = \arg \min_{g_k \in -\mathcal{C}_0^1} \mathbb{E} \left(f(X) - \sum_{k'} f_{k'}^*(X_{k'}) - g_k(X_k) \right)^2$$

We screen out S^C , any variable k that is zero after the second stage, and output S .

$$(3.25) \quad S^C = \{k : f_k^* = 0 \text{ and } g_k^* = 0\}.$$

We refer to this procedure as AC/DC (additive convex/decoupled concave). Theorem 3.2 guarantees that the true set of relevant variables S_0 must be a subset of S .

It is straightforward to construct a finite sample variable screening procedure, which we describe in Figure 3. We use an ℓ_∞/ℓ_1 penalty in equation (3.26) and an ℓ_∞ penalty in equation (3.24) to encourage sparsity. Other penalties can also produce sparse estimates, such as a penalty on the derivative of each of the component functions. The $\|\cdot\|_\infty$ norm is convenient for both theoretical analysis and implementation.

The optimization in (3.26) appears to be infinite dimensional, but it is equivalent to a finite dimensional quadratic program. In the following section, we give the details of this optimization, and show how it can be reformulated to be more computationally efficient.

4. Optimization. We now describe in detail the optimization algorithm for the additive convex regression stage. The second decoupled concave regression stage follows a very similar procedure.

Let $\mathbf{x}_i \in \mathbb{R}^p$ be the covariate, let y_i be the response and let ϵ_i be the mean zero noise. The regression function $f(\cdot)$ we estimate is the sum of univariate functions $f_k(\cdot)$

Input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, regularization parameter λ .

AC Stage: Estimate a sparse additive convex model:

$$(3.26) \quad \hat{f}_1, \dots, \hat{f}_p, \hat{\mu} = \arg \min_{f_1, \dots, f_p \in \mathcal{C}_0^1} \frac{1}{n} \sum_{i=1}^n \left(y_i - \mu - \sum_{k=1}^p f_k(x_{ik}) \right)^2 + \lambda \sum_{k=1}^p \|f_k\|_\infty$$

DC Stage: Estimate concave functions for each k such that $\|\hat{f}_k\|_\infty = 0$:

$$(3.27) \quad \hat{g}_k = \arg \min_{g_k \in \mathcal{C}_0^1} \frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{\mu} - \sum_{k'} \hat{f}_{k'}(x_{ik'}) - g_k(x_{ik}) \right)^2 + \lambda \|g_k\|_\infty$$

Output: Component functions $\{\hat{f}_k\}$ and relevant variables \hat{S} where

$$(3.28) \quad \hat{S}^c = \{k : \|\hat{f}_k\| = 0 \text{ and } \|\hat{g}_k\| = 0\}.$$

FIG 3. The AC/DC algorithm for variable selection in convex regression. The AC stage fits a sparse additive convex regression model, using a quadratic program that imposes an group sparsity penalty for each component function. The DC stage fits decoupled concave functions on the residuals, for each component that is zeroed out in the AC stage.

in each variable dimension and a scalar offset μ . We impose additional constraints that each function $f_k(\cdot)$ is convex, which can be represented by its supporting hyperplanes, i.e.,

$$(4.1) \quad f_{i'k} \geq f_{ik} + \beta_{ik}(x_{i'k} - x_{ik}) \quad \text{for all } i, i' = 1, \dots, n,$$

where $f_{ik} := f_k(x_{ik})$ is the function value and β_{ik} is a subgradient at point x_{ik} . This ostensibly requires $O(n^2p)$ constraints to impose the supporting hyperplane constraints. In fact, only $O(np)$ constraints suffice, since univariate convex functions are characterized by the condition that the subgradient, which is a scalar, must increase monotonically. This observation leads to the optimization

$$(4.2) \quad \begin{aligned} & \min_{\{f_k, \beta_k\}, \mu} \frac{1}{2n} \sum_{i=1}^n \left(y_i - \mu - \sum_{k=1}^p f_{ik} \right)^2 + \lambda \sum_{k=1}^p \|f_k\|_\infty \\ & \text{subject to for all } k = 1, \dots, p: \\ & f_{\pi_k(i+1)k} = f_{\pi_k(i)k} + \beta_{\pi_k(i)k}(x_{\pi_k(i+1)k} - x_{\pi_k(i)k}), \text{ for } i = 1, \dots, n-1 \\ & \sum_{i=1}^n f_{ik} = 0, \\ & \beta_{\pi_k(i+1)k} \geq \beta_{\pi_k(i)k} \text{ for } i = 1, \dots, n-1. \end{aligned}$$

Here f_k denotes the vector $f_k = (f_{1k}, f_{2k}, \dots, f_{nk})^T \in \mathbb{R}^n$ and $\{\pi_k(1), \pi_k(2), \dots, \pi_k(n)\}$

are the indices in the sorted ordering of the values of coordinate k :

$$(4.3) \quad x_{\pi_k(1)k} \leq x_{\pi_k(2)k} \leq \cdots \leq x_{\pi_k(n)k}.$$

We can solve for μ explicitly as $\mu = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$. This follows from the KKT conditions and the constraints $\sum_i f_{ki} = 0$.

The sparse convex additive model optimization in (4.2) is a quadratic program with $O(np)$ variables and $O(np)$ constraints. Directly applying a QP solver for f and β is computationally expensive for relatively large n and p . However, notice that variables in different feature dimensions are only coupled in the squared error term $(y_i - \mu - \sum_{k=1}^p f_{ik})^2$. Hence, we can apply the block coordinate descent method, where in each step we solve the following QP subproblem for $\{f_k, \beta_k\}$ with the other variables fixed. In matrix notation, the optimization is

$$(4.4) \quad \begin{aligned} \min_{f_k, \beta_k, \gamma_k} \quad & \frac{1}{2n} \|r_k - f_k\|_2^2 + \lambda \gamma_k \\ \text{such that} \quad & P_k f_k = \text{diag}(P_k \mathbf{x}_k) \beta_k \\ & D_k \beta_k \leq 0 \\ & -\gamma_k \mathbf{1}_n \leq f_k \leq \gamma_k \mathbf{1}_n \\ & \mathbf{1}_n^\top f_k = 0 \end{aligned}$$

where $\beta_k \in \mathbb{R}^{n-1}$ is the vector $\beta_k = (\beta_{1k}, \dots, \beta_{(n-1)k})^T$, and $r_k \in \mathbb{R}^n$ is the residual vector $r_k = (y_i - \hat{\mu} - \sum_{k' \neq k} f_{ik'})^T$. In addition, $P_k \in \mathbb{R}^{(n-1) \times n}$ is a permutation matrix where the i -th row is all zeros except for the value -1 in position $\pi_k(i)$ and the value 1 in position $\pi_k(i+1)$, and $D_k \in \mathbb{R}^{(n-2) \times (n-1)}$ is another permutation matrix where the i -th row is all zeros except for a value 1 in position $\pi_k(i)$ and a value -1 in position $\pi_k(i+1)$. We denote by $\text{diag}(v)$ the diagonal matrix with diagonal entries v . The extra variable γ_k is introduced to impose the regularization penalty involving the ℓ_∞ norm.

This QP subproblem involves $O(n)$ variables, $O(n)$ constraints and a sparse structure, which can be solved efficiently using optimization packages. In our experiments we use MOSEK (www.mosek.com). We cycle through all covariates k from 1 to p multiple times until convergence. Empirically, we observe that the algorithm converges in only a few cycles. We also implemented an ADMM solver for (4.2) (Boyd *et al.*, 2011), but found that it is not as efficient as this blockwise QP solver.

After optimization, the function estimate for an input vector \mathbf{x} is, according to (4.1),

$$(4.5) \quad \hat{f}(\mathbf{x}) = \sum_{k=1}^p \hat{f}_k(x_k) + \hat{\mu} = \sum_{k=1}^p \max_i \left\{ \hat{f}_{ik} + \hat{\beta}_{ik}(x_k - x_{ik}) \right\} + \hat{\mu}.$$

The univariate concave function estimation required in the DC stage is a straightforward modification of optimization (4.4). It is only necessary to modify the linear inequality constraints so that the subgradients are non-increasing: $\beta_{\pi_k(i+1)k} \leq \beta_{\pi_k(i)k}$.

4.1. *Alternative Formulation.* Optimization (4.2) can be reformulated in terms of the second derivatives. The alternative formulation replaces the order constraints $\beta_{\pi_k(i+1)k} \geq \beta_{\pi_k(i)k}$ with positivity constraints, which simplifies the analysis.

Define $d_{\pi_k(i)k}$ as the second derivative: $d_{\pi_k(1)k} = \beta_{\pi_k(1)k}$, and $d_{\pi_k(i)k} = \beta_{\pi_k(i)k} - \beta_{\pi_k(i-1)k}$ for $i > 1$. The convexity constraint is equivalent to the constraint that $d_{\pi_k(i)k} \geq 0$ for all $i > 1$.

It is easy to verify that $\beta_{\pi_k(i)k} = \sum_{j \leq i} d_{\pi_k(j)k}$ and

$$\begin{aligned} f_k(x_{\pi_k(i)k}) &= f_k(x_{\pi_k(i-1)k}) + \beta_{\pi_k(i-1)k}(x_{\pi_k(i)k} - x_{\pi_k(i-1)k}) \\ &= f_k(x_{\pi_k(1)k}) + \sum_{j < i} \beta_{\pi_k(j)k}(x_{\pi_k(j)k} - x_{\pi_k(j-1)k}) \\ &= f_k(x_{\pi_k(1)k}) + \sum_{j < i} \sum_{j' \leq j} d_{\pi_k(j')k}(x_{\pi_k(j)k} - x_{\pi_k(j-1)k}) \\ &= f_k(x_{\pi_k(1)k}) + \sum_{j' < i} d_{\pi_k(j')k} \sum_{i > j \geq j'} (x_{\pi_k(j)k} - x_{\pi_k(j-1)k}) \\ &= f_k(x_{\pi_k(1)k}) + \sum_{j' < i} d_{\pi_k(j')k}(x_{\pi_k(i)k} - x_{\pi_k(j')k}). \end{aligned}$$

We can write this more compactly in matrix notation as

$$\begin{bmatrix} f_k(x_{\pi_k(1)k}) \\ f_k(x_{\pi_k(2)k}) \\ \vdots \\ f_k(x_{\pi_k(n)k}) \end{bmatrix} = \begin{bmatrix} (x_{1k} - x_{\pi_k(1)k})_+ & \cdots & (x_{1k} - x_{\pi_k(n-1)k})_+ \\ \vdots & \ddots & \vdots \\ (x_{nk} - x_{\pi_k(1)k})_+ & \cdots & (x_{nk} - x_{\pi_k(n-1)k})_+ \end{bmatrix} \begin{bmatrix} d_{\pi_k(1)k} \\ \vdots \\ d_{\pi_k(n-1)k} \end{bmatrix} + \mu_k$$

$$(4.6) \quad \equiv \Delta_k d_k + \mu_k$$

where Δ_k is a $n \times n-1$ matrix such that $\Delta_k(i, j) = (x_{ik} - x_{\pi_k(j)k})_+$, $d_k = (d_{\pi_k(1)k}, \dots, d_{\pi_k(n-1)k})$, and $\mu_k = f_k(x_{\pi_k(1)k})\mathbf{1}_n$. Because f_k has to be centered, $\mu_k = -\frac{1}{n}\mathbf{1}_n^\top \Delta_k d_k$, and therefore

$$(4.7) \quad \Delta_k d_k + \mu_k \mathbf{1}_n = \Delta_k d_k - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \Delta_k d_k = \bar{\Delta}_k d_k$$

where $\bar{\Delta}_k \equiv \Delta_k - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \Delta_k$ is Δ_k with the mean of each column subtracted.

The above derivations prove the following proposition, which states that (4.2) has an alternative formulation.

PROPOSITION 4.1. *Let $\{\hat{f}_k, \hat{\beta}_k\}_{k=1, \dots, p}$ be an optimal solution to (4.2) and suppose $\bar{Y} = 0$. Define vectors $\hat{d}_k \in \mathbb{R}^{n-1}$ such that $\hat{d}_{\pi_k(1)k} = \hat{\beta}_{\pi_k(1)k}$ and $\hat{d}_{\pi_k(i)k} = \hat{\beta}_{\pi_k(i)k} - \hat{\beta}_{\pi_k(i-1)k}$ for $i > 1$. Then $\hat{f}_k = \bar{\Delta}_k \hat{d}_k$ and \hat{d}_k is an optimal solution to the following optimization:*

$$(4.8) \quad \min_{\{d_k \in \mathbb{R}^{n-1}\}_{k=1, \dots, p}} \frac{1}{2n} \left\| Y - \sum_{k=1}^p \bar{\Delta}_k d_k \right\|_2^2 + \lambda_n \sum_{k=1}^p \|\bar{\Delta}_k d_k\|_\infty$$

such that $d_{\pi_k(2)k}, \dots, d_{\pi_k(n-1)k} \geq 0$ (convexity).

Likewise, suppose $\{\hat{d}_k\}_{k=1,\dots,p}$ is a solution to (4.8), define $\hat{\beta}_{\pi_k(i)k} = \sum_{j \leq i} \hat{d}_{\pi_k(j)k}$ and $\hat{f}_k = \bar{\Delta}_k \hat{d}_k$. Then $\{\hat{f}_k, \hat{\beta}_k\}_{k=1,\dots,p}$ is an optimal solution to (4.2). $\bar{\Delta}$ is the n by $n-1$ matrix defined by (4.7).

The decoupled concave postprocessing stage optimization is again similar. Specifically, suppose \hat{d}_k is the output of optimization (4.8), and define the residual vector

$$(4.9) \quad \hat{r} = Y - \sum_{k=1}^p \bar{\Delta}_k \hat{d}_k.$$

Then for all k such that $\hat{d}_k = 0$, the DC stage optimization is formulated as

$$(4.10) \quad \min_{c_k} \frac{1}{2n} \left\| \hat{r} - \Delta_k c_k \right\|_2^2 + \lambda_n \|\Delta_k c_k\|_\infty$$

such that $c_{\pi_k(2)k}, \dots, c_{\pi_k(n-1)k} \leq 0$ (concavity).

We can use either the off-centered Δ_k matrix or the centered $\bar{\Delta}_k$ matrix because the concave estimations are decoupled and hence are not subject to non-identifiability under additive constants.

5. Analysis of Variable Screening Consistency. Our goal is to show that variable screening consistency. That is, as $n, p \rightarrow \infty$, $\mathbb{P}(\hat{S} = S)$ approaches 1 where \hat{S} is the set of variables outputted by AC/DC in the finite sample setting (Algorithm 3) and S is the set of variables outputted in the population setting (3.25).

We divide our analysis into two parts. We first establish a sufficient deterministic condition for consistency of the sparsity pattern screening procedure. We then consider the stochastic setting and argue that the deterministic conditions hold with high probability. Note that in all of our results and analysis, we let c, C represent absolute constants; the actual values of c, C may change from line to line. We derived two equivalent optimizations for AC/DC: (4.2) outputs \hat{f}_k, \hat{g}_k and (4.8) outputs the second derivatives \hat{d}_k . Their equivalence is established in Proposition 4.1 and we use both \hat{d}_k and \hat{f}_k in our analysis. We will also assume in this section that the true regression function f_0 has mean-zero and therefore, we will omit the intercept term \widehat{mu} from our estimation procedure.

In our analysis, we assume that an upper bound B to $\|\hat{f}_k\|_\infty$ is imposed in the optimization procedure, where B is chosen to also upper bound $\|f_k^*\|_\infty$. This B -boundedness constraint is added so that we may use the convex function bracketing results from Kim and Samworth (2014) to establish uniform convergence between the empirical risk and the population risk. We emphasize that this constraint is not needed in practice and we do not use it for any of our simulations.

5.1. Deterministic Setting. We analyze Optimization 4.8 and construct an additive convex solution $\{\hat{d}_k\}_{k=1,\dots,p}$ that is zero for $k \in S^c$, where S is the set of relevant variables, and show that it satisfies the KKT conditions for optimality of

optimization (4.8). We define \hat{d}_k for $k \in S$ to be a solution to the restricted regression (defined below). We also show that $\hat{c}_k = 0$ satisfies the optimality condition of optimization (4.10) for all $k \in S^c$.

DEFINITION 5.1. We define the *restricted regression* problem

$$\min_{d_k} \frac{1}{n} \left\| Y - \sum_{k \in S} \bar{\Delta}_k d_k \right\|_2^2 + \lambda_n \sum_{k \in S} \|\bar{\Delta}_k d_k\|_\infty \quad \text{such that } d_{k,1}, \dots, d_{k,n-1} \geq 0$$

where we restrict the indices k in optimization (4.8) to lie in some set S which contains the true relevant variables.

THEOREM 5.1 (Deterministic setting). *Let $\{\hat{d}_k\}_{k \in S}$ be a minimizer of the restricted regression as defined above. Let $\hat{r} := Y - \sum_{k \in S} \bar{\Delta}_k \hat{d}_k$ be the restricted regression residual.*

Let $\pi_k(i)$ be a reordering of X_k in ascending order so that $X_{k\pi_k(n)}$ is the largest entry. Let $\mathbf{1}_{\pi_k(i:n)}$ be 1 on the coordinates $\pi_k(i), \pi_k(i+1), \dots, \pi_k(n)$ and 0 elsewhere. Define $\text{range}_k = X_{k\pi_k(n)} - X_{k\pi_k(1)}$.

Suppose for all $k \in S^c$, for all $i = 1, \dots, n$, $\lambda_n > \text{range}_k \left| \frac{32}{n} \hat{r}^\top \mathbf{1}_{\pi_k(i:n)} \right|$. Suppose also that for all $k \in S^c$, $\max_{i=1, \dots, n-1} \frac{X_{k\pi_k(i+1)} - X_{k\pi_k(i)}}{\text{range}_k} \geq \frac{1}{16}$, and $\text{range}_k \geq 1$.

Then the following two statements hold.

1. *Let $\hat{d}_k = 0$ for $k \in S^c$. Then $\{\hat{d}_k\}_{k=1, \dots, p}$ is an optimal solution to optimization (4.8). Furthermore, any solution to the optimization program (4.8) must be zero on S^c .*
2. *For all $k \in S^c$, the solution \hat{c}_k to optimization (4.10) must be zero and unique.*

Theorem 5.1 states that the estimator produces no false positive so long as λ_n upper bounds the partial sums of the residual \hat{r} and that the maximum gap between ordered values of X_k is small.

This result holds regardless of whether or not we impose the boundedness conditions in optimization (4.8) and (4.10). The full proof of Theorem 5.1 is in Section 8.1 of the Appendix. We allow S in Theorem 5.1 to be any set containing the relevant variables; in Lasso analysis, S is taken to be the set of relevant variables; we will take S to be the set of variables chosen by the additive convex and decoupled concave procedure in the population setting, which is guaranteed to contain the relevant variables because of additive faithfulness.

Theorem 5.1 allows us to separately analyze the false negative rates and false positive rates. To control false positives, we analyze the condition on λ_n for $k \in S^c$. To control false negatives, we need only analyze the restricted regression with only $|S|$ variables.

The proof of Theorem 5.1 analyses the KKT conditions of optimization (4.8). This parallels the now standard *primal-dual witness* technique (Wainwright, 2009).

However, we cannot derive analogous *mutual incoherence* conditions because the estimation is nonparametric—even the low dimensional restricted regression has $s(n-1)$ variables. The details of the proof are given in Section 8.1 of the Appendix.

5.2. Probabilistic Setting. In the probabilistic setting we treat the covariates as random. We adopt the following standard setup:

1. The data $X^{(1)}, \dots, X^{(n)} \sim F$ are iid from a distribution F with a density $p(\mathbf{x})$ that is supported and strictly positive on $\mathcal{X} = [-1, 1]^p$.
2. The response is $Y = f_0(X) + W$ where W is independent, zero-mean noise; thus $Y^{(i)} = f_0(X^{(i)}) + W^{(i)}$.
3. The regression function f_0 satisfies $f_0(X) = f_0(X_{S_0})$, where $S_0 = \{1, \dots, s_0\}$ is the set of relevant variables.

Let \mathcal{C}^1 denote the set of univariate convex functions supported on $[-1, 1]$, and let \mathcal{C}_1^p denote the set of convex additive functions $\mathcal{C}_1^p \equiv \{f : f = \sum_{k=1}^p f_k, f_k \in \mathcal{C}^1\}$. Let $f^*(\mathbf{x}) = \sum_{k=1}^p f_k^*(x_k)$ be the population risk minimizer in \mathcal{C}_1^p ,

$$(5.1) \quad f^* = \arg \min_{f \in \mathcal{C}_1^p} \mathbb{E}(f_0(X) - f(X))^2.$$

f^* is the unique minimizer by Theorem 3.2. Similarly, we define $-\mathcal{C}^1$ as the set of univariate concave functions supported on $[-1, 1]$ and define

$$(5.2) \quad g_k^* = \arg \min_{g_k \in -\mathcal{C}^1} \mathbb{E}(f_0(X) - f^*(X) - g_k(X_k))^2.$$

The \hat{g}_k 's are unique minimizers as well. We let $S = \{k = 1, \dots, p : f_k^* \neq 0 \text{ or } g_k^* \neq 0\}$ and let $s = |S|$. By additive faithfulness (Theorem 3.2), it must be that $S_0 \subset S$ and thus $s \geq s_0$. In some cases, such as when $X_{S_0}, X_{S_0^c}$ are independent, we have $S = S_0$. Each of our theorems will use a subset of the following assumptions:

A1: X_S, X_{S^c} are independent.

A2: f_0 is convex and twice-differentiable. $\mathbb{E}f_0(X) = 0$.

A3: $\|f_0\|_\infty \leq sB$ and $\|f_k^*\|_\infty \leq B$ for all k .

A4: W is mean-zero sub-Gaussian, independent of X , with scale σ ; i.e., for all $t \in \mathbb{R}$, $\mathbb{E}e^{t\epsilon} \leq e^{\sigma^2 t^2/2}$.

A5: The density $p(\mathbf{x})$ is bounded away from 0 and satisfies the boundary flatness condition. (Definition 3.2)

By assumption A1, f_k^* must be zero for $k \notin S$. We define α_+, α_- as a measure of the signal strength of the weakest variable:

$$(5.3) \quad \begin{aligned} \alpha_+ &= \inf_{f \in \mathcal{C}_1^p : \text{supp}(f) \subsetneq \text{supp}(f^*)} \left\{ \mathbb{E}(f_0(X) - f(X))^2 - \mathbb{E}(f_0(X) - f^*(X))^2 \right\} \\ \alpha_- &= \min_{k \in S : g_k^* \neq 0} \left\{ \mathbb{E}(f_0(X) - f^*(X))^2 - \mathbb{E}(f_0(X) - f^*(X) - g_k^*(X_k))^2 \right\} \end{aligned}$$

α_+ is a lower bound on the excess risk incurred by any additive convex function whose support is strictly smaller than f^* ; $\alpha_+ > 0$ since f^* is the unique risk minimizer. Likewise, α_- lower bounds the excess risk of any decoupled concave fit of the residual $f_0 - f^*$ that is strictly more sparse than the optimal decoupled concave fit $\{\hat{g}_k^*\}$; $\alpha_- > 0$ by the uniqueness of $\{g_k^*\}$ as well. These quantities play the role of the absolute value of the smallest nonzero coefficient in the true linear model in lasso theory. Intuitively, if α_+ is small, then it is easier to make a false omission in the additive convex stage of the procedure. If α_- is small, then it is easier to make a false omission in the decoupled concave stage of the procedure.

REMARK 5.1. We make strong assumptions on the covariates in A1 in order to make weak assumptions on the true regression function f_0 in A2. In particular, we do not assume that f_0 is additive. An important direction for future work is to weaken assumption A1. Our simulation experiments indicate that the procedure can be effective even when the relevant and irrelevant variables are correlated.

THEOREM 5.2 (Controlling false positives). *Suppose assumptions A1-A5 hold. Define $\tilde{\sigma} \equiv \max(\sigma, B)$ and define $\text{range}_k = X_{k\pi_k(n)} - X_{k\pi_k(1)}$. Suppose $p \leq O(\exp(cn))$ and $n \geq C$ for some positive constants $0 < c < 1$ and C . Suppose also*

$$(5.4) \quad \lambda_n \geq 512s\tilde{\sigma}\sqrt{\frac{\log^2 np}{n}}.$$

Then with probability at least $1 - \frac{24}{n}$, for all $k \in S^c$, for all $i = 1, \dots, n$,

$$(5.5) \quad \lambda_n \geq \text{range}_k \left| \frac{32}{n} \hat{r}^\top \mathbf{1}_{(i':n)_k} \right|,$$

$\max_{i'} \frac{X_{k\pi_k(i'+1)} - X_{k\pi_k(i')}}{\text{range}_k} \leq \frac{1}{16}$, $\text{range}_k \geq 1$, and both the AC solution \hat{f}_k from optimization (4.8) and the DC solution \hat{g}_k from optimization (4.10) are zero.

The proof of Theorem 5.2 exploits independence of \hat{r} and X_k under assumption A1; when \hat{r} and X_k are independent, $\hat{r}^\top \mathbf{1}_{(i':n)_k}$ is the sum of $n - i' + 1$ random coordinates of \hat{r} . We can then use concentration of measure results for sampling without replacement to argue that $|\frac{1}{n} \hat{r}^\top \mathbf{1}_{(i':n)_k}|$ is small with high probability. The result of Theorem 5.1 is then used. The full proof of Theorem 5.2 is in Section 8.2 of the Appendix.

THEOREM 5.3 (Controlling false negatives). *Suppose assumptions A1-A5 hold. Let \hat{f} be any AC solution to the restricted regression with B -boundedness constraint, and let \hat{g}_k be any DC solution to the restricted regression with B -boundedness constraint. Let $\tilde{\sigma}$ denote $\max(\sigma, B)$. Suppose*

$$(5.6) \quad \lambda_n \leq 512s\tilde{\sigma}\sqrt{\frac{\log^2 np}{n}}$$

and that n is sufficiently large so that, for some constant $c' > 1$,

$$(5.7) \quad \frac{n^{4/5}}{\log np} \geq c' B^4 \tilde{\sigma}^2 s^5.$$

Assume that the signal-to-noise ratio satisfies

$$(5.8) \quad \frac{\alpha_+}{\tilde{\sigma}} \geq cB^2 \sqrt{\frac{s^5}{n^{4/5}} \log^2 np}$$

$$(5.9) \quad \frac{\alpha_-^2}{\tilde{\sigma}} \geq cB^2 \sqrt{\frac{s^5}{n^{4/5}} \log^2 np}$$

where c is a constant. Then with probability at least $1 - \frac{C}{n}$ for some constant C , $\hat{f}_k \neq 0$ or $\hat{g}_k \neq 0$ for all $k \in S$.

This is a finite sample version of Theorem 3.1. We need stronger assumptions in Theorem 5.3 to use our additive faithfulness result, Theorem 3.1. The full proof of Theorem 5.3 is in Section 8.3 of the appendix.

Combining Theorems 5.2 and 5.3 we obtain the following result.

COROLLARY 5.1. *Suppose the assumptions of Theorem 5.2 and Theorem 5.3 hold. Then with probability at least $1 - \frac{C}{n}$*

$$(5.10) \quad \hat{f}_k \neq 0 \text{ or } \hat{g}_k \neq 0 \text{ for all } k \in S$$

$$(5.11) \quad \hat{f}_k = 0 \text{ and } \hat{g}_k = 0 \text{ for all } k \notin S$$

for some constant C .

The above corollary implies that consistent variable selection is achievable with an exponential scaling of the ambient dimension scaling, $p = O(\exp(n^c))$ for some $0 < c < 1$, just as in parametric models. The cost of nonparametric modeling through shape constraints is reflected in the scaling with respect to the number of relevant variables, which can scale as $s = o(n^{4/25})$.

REMARK 5.2. [Comminges and Dalalyan \(2012\)](#) have shown that under traditional smoothness constraints, even with a product distribution, variable selection is achievable only if $n > O(e^{s_0})$. It is interesting to observe that because of additive faithfulness, the convexity assumption enables a much better scaling of $n = O(\text{poly}(s_0))$, demonstrating that geometric constraints can be quite different from the previously studied smoothness conditions.

6. Experiments. We first illustrate our methods using a simulation of the model

$$Y_i = x_{iS}^\top Q x_{iS} + \epsilon_i \quad (i = 1, 2, \dots, n).$$

Here x_i denotes data sample i drawn from $\mathcal{N}(0, \Sigma)$, x_{iS} is a subset of x_i with dimension $|S| = 5$, where S represents the relevant variable set, and ϵ_i is additive noise drawn from $\mathcal{N}(0, 1)$. The matrix Q is a symmetric positive definite matrix of dimension $|S| \times |S|$. Note that if Q is diagonal, then the true function is convex additive; otherwise the true function is convex but not additive. For all simulations in this section, we set $\lambda = 4\sqrt{\log(np)/n}$.

In the first simulation, we vary the ambient dimension p . We set Q as one on the diagonal and $1/2$ on the off-diagonal with 0.5 probability, and choose $n = 100, 200, \dots, 1000$ and $p = 64, 128, 256$ and 512 . We use independent design by setting $\Sigma = I_p$. For each (n, p) combination, we generate 100 independent data sets. For each data set we use the AC/DC algorithm and mark feature j as irrelevant if both the AC estimate $\|\hat{f}_j\|_\infty$ and the DC estimate $\|\hat{g}_k\|_\infty$ are smaller than 10^{-6} . We plot the probability of exact support recovery over the 100 data sets in Figure 4(a). We observe that the algorithm performs consistent variable selection even if the dimensionality is large. To give the reader a sense of the running speed, for a data set with $n = 1000$ and $p = 512$, the code runs in roughly two minutes on a machine with 2.3 GHz Intel Core i5 CPU and 4 GB memory.

In the second simulation, we vary the sparsity of the Q matrix, that is, we vary the extent to which the true function is non-additive. We generate four Q matrices plotted in Figure 4(b), where the diagonal elements are all one and the off-diagonal elements are $\frac{1}{2}$ with probability α ($\alpha = 0, 0.2, 0.5, 1$ for the four cases). We fix $p = 128$ and choose $n = 100, 200, \dots, 1000$. We use independent design by setting $\Sigma = I_p$. We again run the AC/DC optimization on 100 independently generated data sets and plot the probability of recovery in Figure 4(c). The results demonstrate that AC/DC performs consistent variable selection even if the true function is not additive (but still convex).

In the third simulation, we use correlated design and vary the correlation. We let x_i be drawn from $\mathcal{N}(0, \Sigma)$ instead of $\mathcal{N}(0, I_p)$, with $\Sigma_{ij} = \nu^{|i-j|}$. We use the non-additive Q , same as in the second experiment, with $\alpha = 0.5$ and fix $p = 128$. The recovery curves for $\nu = 0.2, 0.4, 0.6, 0.8$ are shown in Figure 4(d). As can be seen, for design of moderate correlation, AC/DC can still select relevant variables well.

We next use the Boston housing data rather than simulated data. This data set contains 13 covariates, 506 samples and one response variable indicating housing values in suburbs of Boston. The data and detailed description can be found on the UCI Machine Learning Repository website¹.

We first use all $n = 506$ samples (with standardization) in the AC/DC algorithm, using a set of candidate regularization parameters $\{\lambda^{(t)}\}$ ranging from $\lambda^{(1)} = 0$ (no regularization) to 2. For each $\lambda^{(t)}$ we obtain a function value matrix $h^{(t)}$ with $p = 13$ columns. The non-zero columns in this matrix indicate the variables selected using $\lambda^{(t)}$. We plot $\|h^{(t)}\|_\infty$ and the column-wise mean of $h^{(t)}$ versus the normalized norm

¹<http://archive.ics.uci.edu/ml/datasets/Housing>

$\frac{\|h^{(t)}\|_{\infty,1}}{\|h^{(1)}\|_{\infty,1}}$ in Figures 5(a) and 5(b). For comparison, we plot the LASSO/LARS result in a similar way in Figure 5(c). From the figures we observe that the first three variables selected by AC/DC and LASSO are the same: LSTAT, RM and PTRATIO, consistent with previous findings (Ravikumar *et al.*, 2007). The fourth variable selected by AC/DC is INDUS (with $\lambda^{(t)} = 0.7$). We then refit AC/DC with only these four variables without regularization, and plot the estimated additive functions in Figure 5(e). When refitting, we constrain a component to be convex if it is non-zero in the AC stage and concave if it is non-zero in the DC stage. As can be seen, these functions contain clear nonlinear effects which cannot be captured by LASSO. The shapes of these functions, including the concave shape of the PTRATIO function, are in agreement with those obtained by SpAM (Ravikumar *et al.*, 2007).

Next, in order to quantitatively study the predictive performance, we run 3 times 5-fold cross validation, following the same procedure described above—training, variable selection and refitting. A plot of the mean and standard deviation of the predictive mean squared error (MSE) is shown in Figure 5(d). Since for AC/DC the same regularization level $\lambda^{(t)}$ may lead to a slightly different number of selected variables in different folds and runs, the values on the x -axis for AC/DC are not necessarily integers. The figure clearly shows that AC/DC has a lower predictive MSE than LASSO. We also compared the performance of AC/DC with that of Additive Forward Regression (AFR) presented in Liu and Chen (2009), and found that they are similar. The main advantages of AC/DC compared with AFR and SpAM are that there are no smoothing parameters required, and the optimization is formulated as a convex program, guaranteeing a global optimum.

7. Discussion. We have introduced a framework for estimating high dimensional but sparse convex functions. Because of the special properties of convexity, variable selection for convex functions enjoys additive faithfulness—it suffices to carry out variable selection over an additive model, in spite of the approximation error this introduces. Sparse convex additive models can be optimized using block coordinate quadratic programming, which we have found to be effective and scalable. We established variable selection consistency results, allowing exponential scaling in the ambient dimension. We expect that the technical assumptions we have used in these analyses can be weakened; this is one direction for future work. Another interesting direction for building on this work is to allow for additive models that are a combination of convex and concave components. If the convexity/concavity of each component function is known, this again yields a convex program. The challenge is to develop a method to automatically detect the concavity or convexity pattern of the variables.

Acknowledgements. Research supported in part by NSF grants IIS-1116730, AFOSR grant FA9550-09-1-0373, ONR grant N000141210762, and an Amazon AWS in Education Machine Learning Research grant.

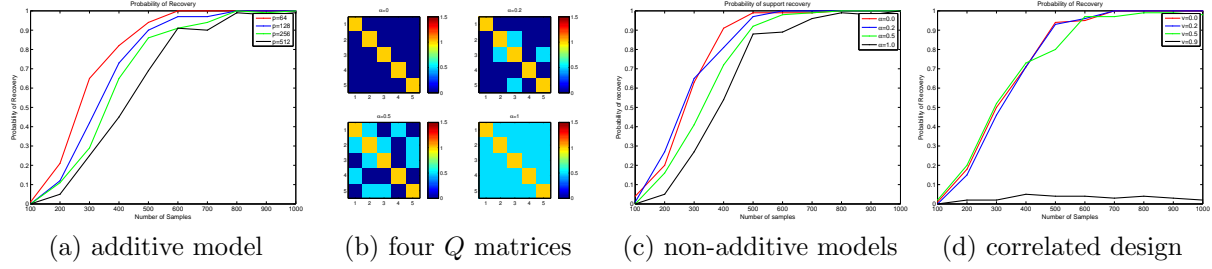


FIG 4. Support recovery results where the additive assumption is correct (a), incorrect (b), (c), and with correlated design (d).

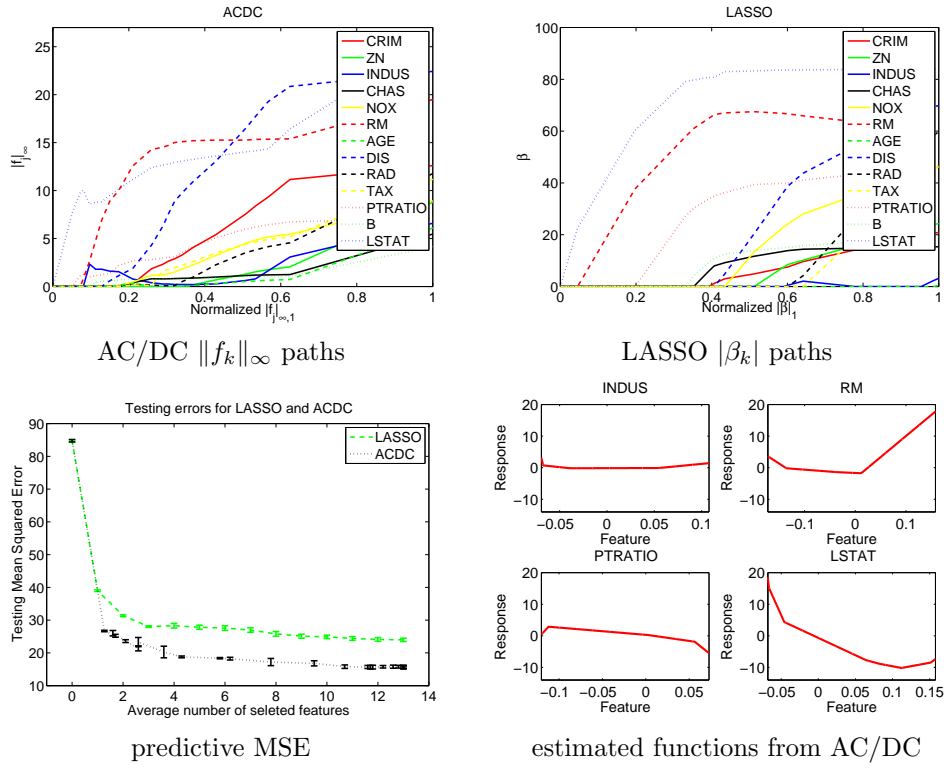


FIG 5. Results on Boston housing data, showing regularization paths, MSE and fitted functions.

References.

- BERTIN, K. and LECUÉ, G. (2008). Selection of variables and dimension reduction in high-dimensional non-parametric regression. *Electronic Journal of Statistics* **2** 1224–1241.
- BOYD, S. and VANDENBERGHE, L. (2004). *Convex Optimization*. Cambridge University Press.
- BOYD, S., PARIKH, N., CHU, E., PELEATO, B. and ECKSTEIN, J. (2011). Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends in Machine Learning* **3** 1–122.
- CHEN, Y. and SAMWORTH, R. J. (2014). Generalised additive and index models with shape constraints. *arXiv preprint arXiv:1404.2957*.
- CHEN, H. and YAO, D. D. (2001). *Fundamentals of Queueing Networks: Performance, Asymptotics, and Optimization*. Springer-Verlag.
- COMMINGES, L. and DALALYAN, A. S. (2012). Tight conditions for consistency of variable selection in the context of high dimensionality. *The Annals of Statistics* **40** 2667–2696.
- CULE, M., SAMWORTH, R. and STEWART, M. (2010). Maximum likelihood estimation of a multi-dimensional log-concave density (with discussion). *J. R. Stat. Soc. Ser. B, Stat. Methodol.* **72** 545–600.
- DEVORE, R., PETROVA, G. and WOJTASZCZYK, P. (2011). Approximation of functions of few variables in high dimensions. *Constructive Approximation* **33** 125–143.
- GOLDENSHLUGER, A. and ZEEVI, A. (2006). Recovering convex boundaries from blurred and noisy observations. *Ann. Statist* **34** 1375–1394.
- GUNTUBOYINA, A. and SEN, B. (2013). Global risk bounds and adaptation in univariate convex regression. *arXiv:1305.1648*.
- HANNAH, L. A. and DUNSON, D. B. (2012). Ensemble Methods for Convex Regression with Applications to Geometric Programming Based Circuit Design. In *International Conference on Machine Learning (ICML)*.
- HORN, R. and JOHNSON, C. (1990). *Matrix Analysis*. Cambridge University Press; Reprint edition.
- HUANG, J., HOROWITZ, J. L. and WEI, F. (2010). Variable selection in nonparametric additive models. *Annals of statistics* **38** 2282.
- KIM, A. K. and SAMWORTH, R. J. (2014). Global rates of convergence in log-concave density estimation. *arXiv preprint arXiv:1404.2298*.
- KOLTCHINSKII, V. and YUAN, M. (2010). Sparsity in multiple kernel learning. *The Annals of Statistics* **38** 3660–3695.
- LAFFERTY, J. and WASSERMAN, L. (2008). Rodeo: sparse, greedy nonparametric regression. *The Annals of Statistics* **36** 28–63.
- LELE, A. S., KULKARNI, S. R. and WILLSKY, A. S. (1992). Convex-polygon estimation from support-line measurements and applications to target reconstruction from laser-radar data. *Journal of the Optical Society of America, Series A* **9** 1693–1714.
- LIM, E. and GLYNN, P. W. (2012). Consistency of Multidimensional Convex Regression. *Operations Research* **60** 196–208.
- LIU, H. and CHEN, X. (2009). Nonparametric greedy algorithm for the sparse learning problems. In *Advances in Neural Information Processing Systems*.
- MEYER, R. F. and PRATT, J. W. (1968). The consistent assessment and fairing of preference functions. *IEEE Trans. Systems Sci. Cybernetics* **4** 270–278.
- MOSSEL, E., O’DONNELL, R. and SERVEDIO, R. (2004). Learning functions of k relevant variables. *Journal of Computer and System Sciences* **69** 421–434.
- PRINCE, J. L. and WILLSKY, A. S. (1990). Reconstructing convex sets from support line measurements. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **12** 377–389.
- PYA, N. and WOOD, S. N. (2014). Shape constrained additive models. *Statistics and Computing* 1–17.
- RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2012). Minimax-optimal rates for sparse additive

- models over kernel classes via convex programming. *J. Mach. Learn. Res.* **13** 389–427.
- RAVIKUMAR, P., LIU, H., LAFFERTY, J. and WASSERMAN, L. (2007). Spam: Sparse additive models. In *Advances in Neural Information Processing Systems*.
- SEIJO, E. and SEN, B. (2011). Nonparametric least squares estimation of a multivariate convex regression function. *The Annals of Statistics* **39** 1633–1657.
- SERFLING, R. J. (1974). Probability inequalities for the sum in sampling without replacement. *The Annals of Statistics* **2** 39–48.
- VERSHYNIN, R. (2010). Introduction to the non-asymptotic analysis of random matrices. arXiv:1011.3027.
- WAINWRIGHT, M. J. (2009). Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (lasso). *Information Theory, IEEE Transactions on* **55** 2183–2202.

8. Supplement: Proofs of Technical Results.

8.1. *Proof of the Deterministic Condition for Sparsistency.* We restate Theorem 5.1 first for convenience. The following holds regardless of whether we impose the B -boundedness condition (see discussion at beginning of Section 5 for definition of the B -boundedness condition).

THEOREM 8.1. *Let $\{\hat{d}_k\}_{k \in S}$ be a minimizer of the restricted regression, that is, the solution to optimization (4.8) where we restrict $k \in S$. Let $\hat{r} := Y - \sum_{k \in S} \bar{\Delta}_k \hat{d}_k$ be the restricted regression residual.*

Let $\pi_k(i)$ be an reordering of X_k in ascending order so that $X_{k\pi_k(n)}$ is the largest entry. Let $\mathbf{1}_{\pi_k(i:n)}$ be 1 on the coordinates $\pi_k(i), \pi_k(i+1), \dots, \pi_k(n)$ and 0 elsewhere. Define $\text{range}_k = X_{k\pi_k(n)} - X_{k\pi_k(1)}$.

Suppose for all $k \in S^c$ and for all $i = 1, \dots, n$, $\lambda_n \geq \text{range}_k \left| \frac{32}{n} \hat{r}^\top \mathbf{1}_{\pi_k(i:n)} \right|$, and $\max_{i=1, \dots, n-1} \frac{X_{k\pi_k(i+1)} - X_{k\pi_k(i)}}{\text{range}_k} \geq \frac{1}{16}$, and $\text{range}_k \geq 1$.

Then the following are true:

1. *Let $\hat{d}_k = 0$ for $k \in S^c$, then $\{\hat{d}_k\}_{k=1, \dots, p}$ is an optimal solution to optimization 4.8. Furthermore, any solution to the optimization program 4.8 must be zero on S^c .*
2. *For all $k \in S^c$, the solution to optimization 4.10 must be zero and unique.*

PROOF. We will omit the B -boundedness constraint in our proof here. It is easy to verify that the result of the theorem still holds with the constraint added in.

We begin by considering the first item in the conclusion of the theorem. We will show that with $\{\hat{d}_k\}_{k=1, \dots, p}$ as constructed, we can set the dual variables to satisfy the complementary slackness and stationary conditions: $\nabla_{d_k} \mathcal{L}(\hat{d}) = 0$ for all k .

The Lagrangian is

$$(8.1) \quad \mathcal{L}(\{d_k\}, \nu) = \frac{1}{2n} \left\| Y - \sum_{k=1}^p \bar{\Delta}_k d_k \right\|_2^2 + \lambda \sum_{k=1}^p \|\bar{\Delta}_k d_k\|_\infty - \sum_{k=1}^p \sum_{i=2}^{n-1} \nu_{ki} d_{ki}$$

with the constraint that $\nu_{ki} \geq 0$ for all k, i .

Because $\{\hat{d}_k\}_{k \in S}$ is by definition the optimal solution of the restricted regression, it is a consequence that stationarity holds for $k \in S$, that is, $\partial_{\{d_k\}_{k \in S}} \mathcal{L}(d) = 0$, and that the dual variables ν_k for $k \in S$ satisfy complementary slackness.

We now verify that stationarity holds also for $k \in S^c$. We fix one dimension $k \in S^c$ and let $\hat{r} = Y - \sum_{k' \in S} \bar{\Delta}_{k'} \hat{d}_{k'}$. The Lagrangian form of the optimization, in terms of just d_k , is

$$\mathcal{L}(d_k, \nu_k) = \frac{1}{2n} \left\| Y - \sum_{k' \in S} \bar{\Delta}_{k'} d_{k'} - \bar{\Delta}_k d_k \right\|_2^2 + \lambda \|\bar{\Delta}_k d_k\|_\infty - \sum_{i=2}^{n-1} \nu_{ki} d_{ki}$$

with the constraint that $\nu_{ki} \geq 0$ for $i = 2, \dots, n-1$.

The derivative of the Lagrangian is

$$\partial_{d_k} \mathcal{L}(d_k) = -\frac{1}{n} \bar{\Delta}_k^\top (Y - \sum_{k' \in S} \bar{\Delta}_{k'} d_{k'} - \bar{\Delta}_k d_k) + \lambda \bar{\Delta}_k^\top \mathbf{u} - \begin{pmatrix} 0 \\ \nu_k \end{pmatrix}$$

where \mathbf{u} is the subgradient of $\|\bar{\Delta}_k d_k\|_\infty$. If $\bar{\Delta}_k d_k = 0$, then \mathbf{u} can be any vector whose L_1 norm is less than or equal to 1. $\nu_k \geq 0$ is a vector of Lagrangian multipliers. ν_{k1} does not exist because d_{k1} is not constrained to be non-negative.

We now substitute in $d_{k'} = \hat{d}_{k'}$ for $k' \in S$, $d_k = 0$ for $k \in S$, and $r = \hat{r}$ and show that the \mathbf{u}, ν_k dual variables can be set in a way to ensure that stationarity:

$$\partial_{d_k} \mathcal{L}(\hat{d}_k) = -\frac{1}{n} \bar{\Delta}_k^\top \hat{r} + \lambda \bar{\Delta}_k^\top \mathbf{u} - \begin{pmatrix} 0 \\ \nu_k \end{pmatrix} = 0$$

where $\|\mathbf{u}\| \leq 1$ and $\nu_k \geq 0$. It clear that to show stationarity, we only need to show that $[-\frac{1}{n} \bar{\Delta}_k^\top \hat{r} + \lambda \bar{\Delta}_k^\top \mathbf{u}]_j = 0$ for $j = 1$ and ≥ 0 for $j = 2, \dots, n-1$.

To ease notational burden, let us reorder the samples in ascending order so that the i -th sample is the i -th smallest sample. We will from here on write X_{ki} to denote $X_{k\pi_k(i)}$.

Define i^* as the largest index such that $\frac{X_{kn} - X_{ki^*}}{X_{kn} - X_{k1}} \geq 1/2$. We will construct $\mathbf{u} = (a - a', 0, \dots, -a, 0, \dots, a')$ where a, a' are positive scalars, where $-a$ lies at the i^* -th coordinate, and where the coordinates of \mathbf{u} correspond to the new sample ordering.

We define

$$\begin{aligned} \kappa &= \frac{1}{\lambda n} [\Delta_k^\top \hat{r}]_1 \\ a' &= \frac{X_{kn} - X_{k1}}{X_{kn} - X_{ki^*}} \kappa + \frac{X_{ki^*} - X_{k1}}{X_{kn} - X_{ki^*}} \frac{1}{8} \\ a &= \frac{X_{kn} - X_{k1}}{X_{kn} - X_{ki^*}} \kappa + \frac{X_{kn} - X_{k1}}{X_{kn} - X_{ki^*}} \frac{1}{8} \end{aligned}$$

and we verify two facts: first that the KKT stationarity is satisfied and second, that $\|\mathbf{u}\|_1 < 1$ with high probability. Our claim is proved immediately by combining these two facts.

Because \hat{r} and \mathbf{u} are both centered vectors, $\bar{\Delta}_k^\top \hat{r} = \Delta_k^\top \hat{r}$ and likewise for \mathbf{u} . Therefore, we need only show that for $j = 1$, $\lambda [\Delta_k^\top \mathbf{u}]_j = [\frac{1}{n} \Delta_k^\top \hat{r}]_j$ and that for $j = 2, \dots, n-1$, $\lambda [\Delta_k^\top \mathbf{u}]_j \geq [\frac{1}{n} \Delta_k^\top \hat{r}]_j$.

With our explicitly defined form of \mathbf{u} , we can characterize

$$(8.2) \quad [\Delta_k^\top \mathbf{u}]_j = \begin{cases} (-a + a')(X_{ki^*} - X_{kj}) + a'(X_{kn} - X_{ki^*}) & \text{if } j \leq i^* \\ a'(X_{kn} - X_{kj}) & \text{if } j \geq i^* \end{cases}$$

It is straightforward to check that $[\lambda \Delta_k^\top \mathbf{u}]_1 = \lambda \kappa = \frac{1}{n} [\Delta_k^\top \hat{\mathbf{r}}]_1$.

To check that other rows of stationarity condition holds, we characterize $[\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}]_j$:

$$\begin{aligned} [\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}]_j &= \frac{1}{n} \sum_{i>j} (X_{ki} - X_{kj}) \hat{r}_i \\ &= \frac{1}{n} \sum_{i>j} \sum_{j<i' \leq i} \text{gap}_{i'} \hat{r}_i \\ &= \frac{1}{n} \sum_{i'>j} \text{gap}_{i'} \sum_{i \geq i'} \hat{r}_i \\ &= \frac{1}{n} \sum_{i'>j} \text{gap}_{i'} \mathbf{1}_{i':n}^\top \hat{\mathbf{r}} \end{aligned}$$

where we denote $\text{gap}_{i'} = X_{ki'} - X_{k(i'-1)}$.

We pause for a second here to give a summary of our proof strategy. We leverage two critical observations: first, any two adjacent coordinates in the vector $\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}$ cannot differ by too much. Second, we defined a, a' such that $-a + a' = -\frac{1}{8}$ so that $\lambda \Delta_k^\top \mathbf{u}$ is a sequence that strictly increases in the first half (for coordinates in $\{1, \dots, i^*\}$) and strictly decreases in the second half.

We know $\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}$ and $\lambda \Delta_k^\top \mathbf{u}$ are equal in the first coordinate. We will show that the second sequence increases faster than the first sequence which will imply that the second sequence is larger than the first in the first half of the coordinates. We will then work similarly but backwards for the second half.

Following our strategy, We first compare $[\lambda \Delta_k^\top \mathbf{u}]_j$ and $[\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}]_j$ for $j = 1, \dots, i^* - 1$.

For $j = 1, \dots, i^* - 1$, we have that

$$\begin{aligned} \lambda [\Delta_k^\top \mathbf{u}]_{j+1} - \lambda [\Delta_k^\top \mathbf{u}]_j &= \lambda(a - a') \text{gap}_{j+1} \\ &\geq -\text{gap}_{j+1} \frac{1}{n} \mathbf{1}_{(j+1):n}^\top \hat{\mathbf{r}} \\ &= [\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}]_{j+1} - [\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}]_j \end{aligned}$$

The inequality follows because $a - a' = \frac{1}{8}$ and thus $\lambda(a - a') \geq \left| \frac{1}{n} \mathbf{1}_{(j+1):n}^\top \hat{\mathbf{r}} \right|$. Therefore, for all $j = 1, \dots, i^*$:

$$[\lambda \Delta_k^\top \mathbf{u}]_j \geq [\frac{1}{n} \Delta_k^\top \hat{\mathbf{r}}]_j$$

For $j \geq i^*$, we start our comparison from $j = n - 1$. First, we claim that $a' > \frac{1}{32}$. To prove this claim, note that

$$(8.3) \quad |\kappa| = \left| \frac{1}{\lambda n} \sum_{i'>1} \text{gap}_{i'} \mathbf{1}_{i':n}^\top \hat{\mathbf{r}} \right| \leq \frac{1}{X_{kn} - X_{k1}} \frac{1}{32} \sum_{i'>1} \text{gap}_{i'} = \frac{1}{32}$$

and that

$$\frac{X_{kn} - X_{ki^*}}{X_{kn} - X_{k1}} = \frac{X_{kn} - X_{k(i^*+1)} + X_{k(i^*+1)} - X_{ki^*}}{X_{kn} - X_{k1}} \leq \frac{1}{2} + \frac{1}{16}$$

where the inequality follows because we had assumed that $\frac{X_{k(i+1)} - X_{ki}}{X_{k(n)} - X_{k(1)}} \leq \frac{1}{16}$ for all $i = 1, \dots, n-1$.

So, we have

$$\begin{aligned} a' &= \frac{X_{kn} - X_{k1}}{X_{kn} - X_{ki^*}} \kappa + \frac{X_{ki^*} - X_{k1}}{X_{kn} - X_{ki^*}} \frac{1}{8} \\ &= \frac{X_{kn} - X_{k1}}{X_{kn} - X_{ki^*}} \left(\kappa + \frac{X_{ki^*} - X_{k1}}{X_{kn} - X_{k1}} \frac{1}{8} \right) \\ &\geq \frac{X_{kn} - X_{k1}}{X_{kn} - X_{ki^*}} \left(-\frac{1}{32} + \left(\frac{1}{2} - \frac{1}{16} \right) \frac{1}{8} \right) \\ &\geq \frac{1}{1/2 + 1/16} \left(-\frac{1}{32} + \left(\frac{1}{2} - \frac{1}{16} \right) \frac{1}{8} \right) \\ &\geq \frac{1}{32} \end{aligned}$$

In the first inequality of the above derivation, we used the fact that $\frac{X_{ki^*} - X_{k1}}{X_{kn} - X_{k1}} \leq \frac{1}{2} - \frac{1}{16}$. In the second inequality, we used the fact that the quantity inside the parentheses is positive and $\frac{X_{kn} - X_{k1}}{X_{kn} - X_{ki^*}} \geq \frac{1}{1/2 + 1/16}$.

Now consider $j = n-1$.

$$\left[\frac{1}{n} \Delta_k^\top \hat{r} \right]_{n-1} = \frac{1}{n} \text{gap}_n \hat{r}_n \leq \text{gap}_n \frac{\lambda}{32} \leq \lambda \text{gap}_n a' = \lambda [\Delta_k^\top \mathbf{u}]_{n-1}$$

For $j = i^*, \dots, n-2$, we have that

$$\begin{aligned} \lambda [\Delta_k^\top \mathbf{u}]_j - \lambda [\Delta_k^\top \mathbf{u}]_{j+1} &= \lambda a' \text{gap}_{j+1} \\ &\geq \text{gap}_{j+1} \frac{1}{n} \mathbf{1}_{(j+1):n}^\top \hat{r} \\ &\geq \left[\frac{1}{n} \Delta_k^\top \hat{r} \right]_j - \left[\frac{1}{n} \Delta_k^\top \hat{r} \right]_{j+1} \end{aligned}$$

Therefore, for $j = i^*, \dots, n-2$,

$$\lambda [\Delta_k^\top \mathbf{u}]_j \geq \frac{1}{n} [\Delta_k^\top \hat{r}]_j$$

We conclude then that $\lambda [\Delta_k^\top \mathbf{u}]_j \geq \left[\frac{1}{n} \Delta_k^\top \hat{r} \right]_j$ for all $j = 2, \dots, n-1$.

We have thus verified that the stationarity equations hold and now will bound $\|\mathbf{u}\|_1$.

$$\|\mathbf{u}\|_1 = |a - a'| + a + a' \leq \frac{1}{8} + 2a \leq \frac{1}{8} + 4|\kappa| + \frac{1}{2} \leq \frac{1}{8} + \frac{1}{8} + \frac{1}{2} < 1$$

In the third inequality, we used the fact that $|\kappa| \leq \frac{1}{32}$.

We have thus proven that there exists one solution $\{\hat{d}_k\}_{k=1,\dots,p}$ such that $\hat{d}_k = 0$ for all $k \in S^c$. Furthermore, we have shown that the subgradient variables \mathbf{u}_k of the solution $\{\hat{d}_k\}$ can be chosen such that $\|\mathbf{u}_k\|_1 < 1$ for all $k \in S^c$.

We now prove that if $\{\hat{d}'_k\}_{k=1,\dots,p}$ is another solution, then it must be that $\hat{d}'_k = 0$ for all $k \in S^c$ as well. We first claim that $\sum_{k=1}^p \bar{\Delta}_k \hat{d}_k = \sum_{k=1}^p \bar{\Delta}_k \hat{d}'_k$. If this were not true, then a convex combination of \hat{d}_k, \hat{d}'_k would achieve a strictly lower objective on the quadratic term. More precisely, let $\zeta \in [0, 1]$. If $\sum_{k=1}^p \bar{\Delta}_k \hat{d}'_k \neq \sum_{k=1}^p \bar{\Delta}_k \hat{d}_k$, then $\|Y - \sum_{k=1}^p \bar{\Delta}_k (\hat{d}_k + \zeta(\hat{d}'_k - \hat{d}_k))\|_2^2$ is strongly convex as a function of ν . Thus, it cannot be that \hat{d}_k and \hat{d}'_k both achieve optimal objective, and we have reached a contradiction.

Now, we look at the stationarity condition for both $\{\hat{d}_k\}$ and $\{\hat{d}'_k\}$. Let $\mathbf{u}_k \in \partial\|\bar{\Delta}_k \hat{d}_k\|_\infty$ and let $\mathbf{u}'_k \in \partial\|\bar{\Delta}_k \hat{d}'_k\|_\infty$ be the two sets of subgradients. Let $\{\nu_{ki}\}$ and $\{\nu'_{ki}\}$ be the two sets of positivity dual variables, for $k = 1, \dots, p$ and $i = 1, \dots, n-1$. Note that since there is no positivity constraint on d_{k1} , we let $\nu_{k1} = 0$ always.

Let us define $\bar{\Delta}$, a $n \times p(n-1)$ matrix, to denote the column-wise concatenation of $\{\bar{\Delta}_k\}_k$ and \hat{d} , a $p(n-1)$ dimensional vector, to denote the concatenation of $\{\hat{d}_k\}_k$. With this notation, we can express $\sum_{k=1}^p \bar{\Delta}_k \hat{d}_k = \bar{\Delta} \hat{d}$.

Since both solutions $(\hat{d}, \mathbf{u}, \nu)$ and $(\hat{d}', \mathbf{u}', \nu')$ must satisfy the stationarity condition, we have that

$$\bar{\Delta}^\top (Y - \bar{\Delta} \hat{d}) + \lambda \sum_{k=1}^p \bar{\Delta}_k^\top \mathbf{u}_k - \nu = \bar{\Delta}^\top (Y - \bar{\Delta} \hat{d}') + \lambda \sum_{k=1}^p \bar{\Delta}_k^\top \mathbf{u}'_k - \nu' = 0.$$

Multiplying both sides of the above equation by \hat{d}' ,

$$\hat{d}'^\top \bar{\Delta}^\top (Y - \bar{\Delta} \hat{d}) + \lambda \sum_{k=1}^p \hat{d}'^\top \bar{\Delta}_k^\top \mathbf{u}_k - \hat{d}'^\top \nu = \hat{d}'^\top \bar{\Delta}^\top (Y - \bar{\Delta} \hat{d}') + \lambda \sum_{k=1}^p \hat{d}'^\top \bar{\Delta}_k^\top \mathbf{u}'_k - \hat{d}'^\top \nu'.$$

Since $\bar{\Delta} \hat{d}_k = \bar{\Delta} \hat{d}$, $\hat{d}'^\top \nu' = 0$ (complementary slackness), and $\hat{d}'^\top \bar{\Delta}_k^\top \mathbf{u}'_k = \|\hat{f}'_k\|_\infty$ (where $\hat{f}'_k = \bar{\Delta}_k \hat{d}'_k$), we have that

$$\lambda \sum_{k=1}^p \hat{d}'^\top \bar{\Delta}_k^\top \mathbf{u}_k - \hat{d}'^\top \nu = \lambda \sum_{k=1}^p \|\hat{f}'_k\|_\infty.$$

On one hand, \hat{d}' is a feasible solution so $\hat{d}'^\top \nu \geq 0$ and so

$$\sum_{k=1}^p \hat{d}'^\top \bar{\Delta}_k^\top \mathbf{u}_k \geq \sum_{k=1}^p \|\hat{f}'_k\|_\infty.$$

On the other hand, by Hölder's inequality,

$$\sum_{k=1}^p \hat{d}'^\top \bar{\Delta}_k^\top \mathbf{u}_k \leq \sum_{k=1}^p \|\hat{f}'_k\|_\infty \|\mathbf{u}_k\|_1.$$

Since \mathbf{u}_k can be chosen so that $\|\mathbf{u}_k\|_1 < 1$ for all $k \in S^c$, we would get a contradiction if $\|\hat{f}'_k\|_\infty > 0$ for some $k \in S^c$. We thus conclude that \hat{d}' must follow the same sparsity pattern.

The second item in the theorem concerning optimization 4.10 is proven in exactly the same way. The Lagrangian of optimization 4.10 is

$$\mathcal{L}_{\text{cave}}(d_k, \nu_k) = \frac{1}{2n} \|\hat{r} - \bar{\Delta}_k d_k\|_2^2 + \lambda \|\bar{\Delta}_k d_k\|_\infty + \sum_{k=1}^p \sum_{i=2}^{n-1} \nu_{ki} d_{ki}.$$

with $\nu_{ki} \geq 0$. The same reasoning applies to show that $\hat{d}_k = 0$ satisfies KKT conditions sufficient for optimality. \square

8.2. Proof of False Positive Control. We note that in the following analysis the symbols c, C represent absolute constants. We will often abuse notation and “absorb” new absolute constants into c, C ; the actual value of c, C could thus vary from line to line. We first restate the theorem for convenience.

THEOREM 8.2. *Suppose assumptions A1-A5 hold. Define $\tilde{\sigma} \equiv \max(\sigma, B)$. Suppose that $p \leq O(\exp(cn))$ and $n \geq C$ for some constants $0 < c < 1$ and C . Define $\text{range}_k = X_{k\pi_k(n)} - X_{k\pi_k(1)}$.*

If $\lambda_n \geq 2(8 \cdot 32)s\tilde{\sigma}\sqrt{\frac{1}{n}\log^2 np}$ then, with probability at least $1 - \frac{24}{n}$, for all $k \in S^c$, and for all $i' = 1, \dots, n$

$$\lambda_n > \text{range}_k \left| \frac{32}{n} \hat{r}^\top \mathbf{1}_{\pi_k(i':n)} \right|$$

and $\max_{i'} \frac{X_{k\pi_k(i'+1)} - X_{k\pi_k(i')}}{\text{range}_k} \leq \frac{1}{16}$ and $\text{range}_k \geq 1$.

Therefore, for all $k \in S^c$, both the AC solution \hat{f}_k from optimization 4.8, and the DC solution \hat{g}_k from optimization 4.10 are zero.

PROOF. The key is to note that \hat{r} and $\Delta_{k,j}$ are independent for all $k \in S^c, j = 1, \dots, n$ because \hat{r} is only dependent on X_S .

Fix j and i . Then $\hat{r}^\top \mathbf{1}_{\pi_k(i':n)}$ is the sum of $n - i' + 1$ random coordinates of \hat{r} . We will use Serfling’s theorem on the concentration of measure of sampling without replacement (Corollary 8.2). We must first bound $\|\hat{r}\|_\infty$ and $\frac{1}{n} \sum_{i=1}^n \hat{r}_i$ before we can use Serfling’s results however.

Step 1: Bounding $\|\hat{r}\|_\infty$. We have $\hat{r}_i = f_0(x_i) + w_i - \hat{f}(x_i)$ where $\hat{f}(x_i) = \sum_{k \in S} \bar{\Delta}_k \hat{d}_k$ is the convex additive function outputted by the restricted regression. Note that both $f_0(x_i)$ and $\hat{f}(x_i)$ are bounded by $2sB$. Because w_i is sub-Gaussian, $|w_i| \leq \sigma \sqrt{2 \log \frac{2}{\delta}}$ with probability at least $1 - \delta$. By union bound across $i = 1, \dots, n$, we have that $\|w\|_\infty \leq \sigma \sqrt{2 \log \frac{2n}{\delta}}$ with probability at least $1 - \delta$.

Putting these observations together,

$$(8.4) \quad \begin{aligned} \|\hat{r}\|_\infty &\leq 2sB + \sigma\sqrt{2\log\frac{2n}{\delta}} \\ &\leq 4s\tilde{\sigma}\sqrt{\log\frac{2n}{\delta}} \end{aligned}$$

with probability at least $1 - \delta$, where we have defined $\tilde{\sigma} = \max(\sigma, B)$, and assumed that $\sqrt{\log\frac{2np}{\delta}} \geq 1$. This assumption holds under the conditions in the theorem which state that $p \leq \exp(cn)$ and $n \geq C$ for some small constant c and large constant C .

Step 2: *Bounding $|\frac{1}{n}\hat{r}^\top \mathbf{1}|$.* We have that

$$\begin{aligned} \frac{1}{n}\hat{r}^\top \mathbf{1} &= \frac{1}{n} \sum_{i=1}^n f_0(x_i) + w_i - \hat{f}(x_i) \\ &= \frac{1}{n} \sum_{i=1}^n f_0(x_i) + w_i \quad (\hat{f} \text{ is centered}). \end{aligned}$$

Since $|f_0(x_i)| \leq sB$, the first term $|\frac{1}{n} \sum_{i=1}^n f_0(x_i)|$ is at most $sB\sqrt{\frac{2}{n} \log \frac{2}{\delta}}$ with probability at most $1 - \delta$ by Hoeffding's inequality. Since w_i is sub-Gaussian, the second term $|\frac{1}{n} \sum_{i=1}^n w_i|$ is at most $\sigma\sqrt{\frac{2}{n} \log \frac{2}{\delta}}$ with probability at most $1 - \delta$. Taking a union bound, we have that

$$(8.5) \quad \begin{aligned} \left| \frac{1}{n}\hat{r}^\top \mathbf{1} \right| &\leq sB\sqrt{\frac{2}{n} \log \frac{4}{\delta}} + \sigma\sqrt{\frac{2}{n} \log \frac{4}{\delta}} \\ &\leq 4s\tilde{\sigma}\sqrt{\frac{1}{n} \log \frac{4}{\delta}} \end{aligned}$$

with probability at least $1 - \delta$.

Step 3: *Apply Serfling's theorem.* For any $k \in S^c$, Serfling's theorem states that with probability at least $1 - \delta$

$$\left| \frac{1}{n}\hat{r}^\top \mathbf{1}_{\pi_k(i':n)} \right| \leq 2\|\hat{r}\|_\infty \sqrt{\frac{1}{n} \log \frac{2}{\delta}} + \left| \frac{1}{n}\hat{r}^\top \mathbf{1} \right|$$

We need Serfling's theorem to hold for all $k = 1, \dots, p$ and $i' = 1, \dots, n$. We also need the events that $\|\hat{r}\|_\infty$ and $|\frac{1}{n}\hat{r}^\top \mathbf{1}|$ are small to hold. Using a union bound, with probability at least $1 - \delta$, for all k, i' ,

$$\begin{aligned} \left| \frac{1}{n}\hat{r}^\top \mathbf{1}_{\pi_k(i':n)} \right| &\leq 2\|\hat{r}\|_\infty \sqrt{\frac{1}{n} \log \frac{6np}{\delta}} + \left| \frac{1}{n}\hat{r}^\top \mathbf{1} \right| \\ &\leq 4s\tilde{\sigma}\sqrt{\log \frac{6n}{\delta}} \sqrt{\frac{1}{n} \log \frac{6np}{\delta}} + 4s\tilde{\sigma}\sqrt{\frac{1}{n} \log \frac{12}{\delta}} \\ &\leq 8s\tilde{\sigma}\sqrt{\frac{1}{n} \log^2 \frac{12np}{\delta}} \end{aligned}$$

In the second inequality, we used equation (8.4) and equation (8.5) from steps 1 and 2 respectively. Setting $\delta = \frac{12}{n}$ gives the desired expression.

Finally, we note that $2 \geq (X_{k\pi_k(n)} - X_{k\pi_k(1)})$ since $X_k \subset [-1, 1]$. This concludes the proof for the first part of the theorem.

To prove the second and the third claims, let the interval $[-1, 1]$ be divided into 64 non-overlapping segments each of length $1/32$. Because X_k is drawn from a density with a lower bound $c_l > 0$, the probability that every segment contains some samples X_{ki} 's is at least $1 - 64 \left(1 - \frac{1}{32}c_l\right)^n$. Let \mathcal{E}_k denote the event that every segment contains some samples.

Define $\text{gap}_i = X_{k\pi_k(i+1)} - X_{k\pi_k(i)}$ for $i = 1, \dots, n-1$ and define $\text{gap}_0 = X_{k\pi_k(1)} - (-1)$ and $\text{gap}_n = 1 - X_{k\pi_k(n)}$.

If any $\text{gap}_i \geq \frac{1}{16}$, then gap_i has to contain one of the segments. Therefore, under event \mathcal{E}_k , it must be that $\text{gap}_i \leq \frac{1}{16}$ for all i .

Thus, we have that $\text{range}_k \geq 2 - 1/8 \geq 1$ and that for all i ,

$$\frac{X_{k\pi_k(i+1)} - X_{k\pi_k(i)}}{\text{range}_k} \geq \frac{1/16}{2 - 1/8} \geq 1/16$$

Taking an union bound for each $k \in S^c$, the probability of that all \mathcal{E}_k hold is at least $1 - p64 \left(1 - \frac{1}{32}c_l\right)^n$.

$p64 \left(1 - \frac{1}{32}c_l\right)^n = 64p \exp(-c'n)$ for some positive constants $0 < c' < 1$ dependent on c_l . Therefore, if $p \leq \exp(-cn)$ for some $0 < c < c'$ and if n is larger than some constant C , $64p \exp(-c'n) \leq 64 \exp(-(c' - c'')n) \leq \frac{12}{n}$.

Taking an union bound with the event that λ_n upper bounds the partial sums of \hat{r} and we establish the claim. □

8.3. Proof of False Negative Control. We begin by introducing some notation.

8.3.1. Notation. If $f : \mathbb{R}^s \rightarrow \mathbb{R}$, we define $\|f\|_P \equiv \mathbb{E}f(X)^2$. Given samples X_1, \dots, X_n , we denote $\|f\|_n \equiv \frac{1}{n} \sum_{i=1}^n f(X_i)^2$ and $\langle f, g \rangle_n \equiv \frac{1}{n} \sum_{i=1}^n f(X_i)g(X_i)$.

Let \mathcal{C}^1 denote the set of univariate convex functions supported on $[-1, 1]$. Let $\mathcal{C}_B^1 \equiv \{f \in \mathcal{C}^1 : \|f\|_\infty \leq B\}$ denote the set of B -bounded univariate convex functions. Define \mathcal{C}^s as the set of convex additive functions and \mathcal{C}_B^s likewise as the set of convex additive functions whose components are B -bounded:

$$\begin{aligned} \mathcal{C}^s &\equiv \{f : f = \sum_{k=1}^s f_k, f_k \in \mathcal{C}^1\} \\ \mathcal{C}_B^s &\equiv \{f \in \mathcal{C}^s : f = \sum_{k=1}^s f_k, \|f_k\|_\infty \leq B\}. \end{aligned}$$

Let $f^*(x) = \sum_{k=1}^s f_k^*(x_k)$ be the population risk minimizer:

$$f^* = \arg \min_{f \in \mathcal{C}^s} \|f_0 - f^*\|_P^2$$

We let sB be an upper bound on $\|f_0\|_\infty$ and B be an upper bound on $\|f_k^*\|_\infty$. It follows that $\|\hat{f}^*\|_\infty \leq sB$.

We define \hat{f} as the empirical risk minimizer:

$$\hat{f} = \arg \min \left\{ \|y - f\|_n^2 + \lambda \sum_{k=1}^s \|f_k\|_\infty : f \in \mathcal{C}_B^s, \mathbf{1}_n^\top f_k = 0 \right\}$$

For $k \in \{1, \dots, s\}$, define g_k^* to be the decoupled concave population risk minimizer

$$g_k^* \equiv \arg \min_{g_k \in \mathcal{C}^1} \|f_0 - f^* - g_k\|_P^2.$$

In our proof, we will analyze g_k^* for each k such that $f_k^* = 0$. Likewise, we define the empirical version:

$$\hat{g}_k \equiv \arg \min \left\{ \|f_0 - \hat{f} - g_k\|_n^2 : g_k \in \mathcal{C}_B^1, \mathbf{1}_n^\top g_k = 0 \right\}.$$

By the definition of the AC/DC procedure, \hat{g}_k is defined only for an index k that has zero as the convex additive approximation.

8.3.2. Proof. By additive faithfulness of the AC/DC procedure, it is known that $f_k^* \neq 0$ or $g_k^* \neq 0$ for all $k \in S$. Our argument will be to show that the risk of the AC/DC estimators \hat{f}, \hat{g} tends to the risk of the population optimal functions f^*, g^* :

$$\begin{aligned} \|f_0 - \hat{f}\|_P^2 &= \|f_0 - f^*\|_P^2 + \text{err}_+(n) \\ \|f_0 - f^* - \hat{g}_k\|_P^2 &= \|f_0 - f^* - g_k^*\|_P^2 + \text{err}_-(n) \quad \text{for all } k \in S \text{ where } f_k^* = 0, \end{aligned}$$

where the estimation errors $\text{err}_+(n)$ and $\text{err}_-(n)$ decrease with n at some rate.

Assuming this, suppose that $\hat{f}_k = 0$ and $f_k^* \neq 0$. Then when n is large enough such that $\text{err}_+(n)$ and $\text{err}_-(n)$ are smaller than α_+ and α_- defined in equation (5.3), we reach a contradiction. This is because the risk $\|f_0 - f^*\|_P$ of f^* is strictly larger by α_+ than the risk of the best approximation whose k -th component is constrained to be zero. Similarly, suppose $f_k^* = 0$ and $g_k^* \neq 0$. Then when n is large enough, \hat{g}_k must not be zero.

Theorem 8.3 and Theorem 8.4 characterize $\text{err}_+(n)$ and $\text{err}_-(n)$ respectively.

THEOREM 8.3. *Let $\tilde{\sigma} \equiv \max(\sigma, B)$, and let \hat{f} be the minimizer of the restricted regression with $\lambda \leq 512s\tilde{\sigma}\sqrt{\frac{1}{n}\log^2 np}$. Suppose $n \geq c_1 s\sqrt{sB}$. Then with probability at least $1 - \frac{C}{n}$,*

$$(8.6) \quad \|f_0 - \hat{f}\|_P^2 - \|f_0 - f^*\|_P^2 \leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log^2 np},$$

where c_1 is an absolute constant and c, C are constants possibly dependent on b .

PROOF. Our proof proceeds in three steps. First, we bound the difference of empirical risks $\|f_0 - \hat{f}\|_n^2 - \|f_0 - f^*\|_n^2$. Second, we bound the cross-term in the bound using a bracketing entropy argument for convex function classes. Finally, we combine the previous two steps to complete the argument.

Step 1. The function \hat{f} minimizes the penalized empirical risk by definition. We would thus like to say that the penalized empirical risk of \hat{f} is no larger than that of f^* . We cannot do a direct comparison, however, because the empirical mean $\frac{1}{n} \sum_i f_k^*(x_{ik})$ is close to, but not exactly zero. We thus have to work first with the function $f^* - \bar{f}^*$. We have that

$$\|y - \hat{f}\|_n^2 + \lambda \sum_{k=1}^s \|\hat{f}_k\|_\infty \leq \|y - f^* + \bar{f}^*\|_n^2 + \lambda \sum_{k=1}^s \|f_k^* - \bar{f}_k^*\|_\infty$$

Plugging in $y = f_0 + w$, we obtain

$$\begin{aligned} \|f_0 + w - \hat{f}\|_n^2 + \lambda \sum_{k=1}^s \left(\|\hat{f}_k\|_\infty - \|f_k^* - \bar{f}_k^*\|_\infty \right) &\leq \|f_0 + w - f^* + \bar{f}^*\|_n^2 \\ \|f_0 - \hat{f}\|_n^2 + 2\langle w, f_0 - \hat{f} \rangle_n + \lambda \sum_{k=1}^s \left(\|\hat{f}_k\|_\infty - \|f_k^* - \bar{f}_k^*\|_\infty \right) \\ &\leq \|f_0 - f^* + \bar{f}^*\|_n^2 + 2\langle w, f_0 - f^* + \bar{f}^* \rangle \\ \|f_0 - \hat{f}\|_n^2 - \|f_0 - f^* + \bar{f}^*\|_n^2 + \lambda \sum_{k=1}^s \left(\|\hat{f}_k\|_\infty - \|f_k^* - \bar{f}_k^*\|_\infty \right) &\leq 2\langle w, \hat{f} - f^* + \bar{f}^* \rangle. \end{aligned}$$

The middle term can be bounded under the assumption that $\|f_k^* - \bar{f}_k^*\|_\infty \leq 2B$; thus,

$$\|f_0 - \hat{f}\|_n^2 - \|f_0 - f^* + \bar{f}^*\|_n^2 \leq 2\langle w, \hat{f} - f^* + \bar{f}^* \rangle + \lambda 2sB.$$

Using Lemma 8.2, we can remove \bar{f}^* from the lefthand side. Thus with probability at least $1 - \delta$,

$$(8.7) \quad \|f_0 - \hat{f}\|_n^2 - \|f_0 - f^*\|_n^2 \leq 2\langle w, \hat{f} - f^* + \bar{f}^* \rangle + \lambda 2sB + c(sB)^2 \frac{1}{n} \log \frac{2}{\delta}.$$

Step 2. We now upper bound the cross term $2\langle w, \hat{f} - f^* + \bar{f}^* \rangle$ using bracketing entropy.

Define $\mathcal{G} = \{f - f^* + \bar{f}^* : f \in \mathcal{C}_B^s\}$ as the set of convex additive functions centered around the function $f^* - \bar{f}^*$. By Corollary 8.3, there is an ϵ -bracketing of \mathcal{G} whose size is bounded by $\log N_{[]} (2\epsilon, \mathcal{G}, L_1(P)) \leq sK^{**} \left(\frac{2sB}{\epsilon} \right)^{1/2}$, for all $\epsilon \in (0, sB\epsilon_3]$. Let us suppose condition 8.11 holds. Then, by Corollary 8.4, with probability at least $1 - \delta$, each bracketing pair (h_U, h_L) is close in $L_1(P_n)$ norm, i.e., for all (h_U, h_L) , $\frac{1}{n} \sum_{i=1}^n |h_U(X_i) - h_L(X_i)| \leq 2\epsilon + sB \sqrt{\frac{sK^{**}(2sB)^{1/2} \log \frac{1}{\delta}}{2\epsilon^{1/2}n}}$. We verify at the end of the proof that condition 8.11 indeed holds.

For each $h \in \mathcal{G}$, there exists a pair (h_U, h_L) such that $h_U(X_i) - h_L(X_i) \geq h(X_i) - h_L(X_i) \geq 0$. Therefore, with probability at least $1 - \delta$, uniformly for all $h \in \mathcal{G}$:

$$\frac{1}{n} \sum_{i=1}^n |h(X_i) - h_L(X_i)| \leq \frac{1}{n} \sum_{i=1}^n |h_U(X_i) - h_L(X_i)| \leq 2\epsilon + (sB) \sqrt{\frac{sK^{**}(2sB)^{1/2} \log \frac{1}{\delta}}{2\epsilon^{1/2}n}}.$$

We denote $\epsilon_{n,\delta} \equiv (sB) \sqrt{\frac{sK^{**}(2sB)^{1/2} \log \frac{1}{\delta}}{2\epsilon^{1/2}n}}$. Let $\mathcal{E}_{[]}$ denote the event that for each $h \in \mathcal{G}$, there exists h_L in the ϵ -bracketing such that $\|h - h_L\|_{L_{P_n}} \leq 2\epsilon + \epsilon_{n,\delta}$. Then $\mathcal{E}_{[]}$ has probability at most $1 - \delta$ as shown.

Let $\mathcal{E}_{\|w\|_\infty}$ denote the event that $\|w\|_\infty \leq \sigma \sqrt{2 \log \frac{2n}{\delta}}$. Then $\mathcal{E}_{\|w\|_\infty}$ has probability at most $1 - \delta$. We now take an union bound over $\mathcal{E}_{\|w\|_\infty}$ and $\mathcal{E}_{[]}$ and get that, with probability at most $1 - 2\delta$, for all h

$$|\langle w, h - h_L \rangle_n| \leq \|w\|_\infty \frac{1}{n} \sum_{i=1}^n |h(X_i) - h_L(X_i)| \leq \sigma \sqrt{2 \log \frac{4n}{\delta}} (\epsilon + \epsilon_{n,2\delta}).$$

Because w is a sub-Gaussian random variable, we have that for any fixed vector $h_L = (h_L(X_1), \dots, h_L(X_n))$, with probability at least $1 - \delta$, $|\langle w, h_L \rangle_n| \leq \|h_L\|_n \sigma \sqrt{\frac{2}{n} \log \frac{2}{\delta}}$. Using another union bound, we have that the event $\sup_{h_L} |\langle w, h_L \rangle| \leq sB\sigma \sqrt{\frac{2}{n} \log \frac{2N_{[]}^{2N_{[]}}}{\delta}}$ has probability at most $1 - \delta$.

Putting this together, we have that

$$\begin{aligned} |\langle w, h \rangle_n| &\leq |\langle w, h_L \rangle_n| + |\langle w, h - h_L \rangle_n| \\ |\sup_{h \in \mathcal{G}} \langle w, h \rangle_n| &\leq |\sup_{h_L} \langle w, h_L \rangle_n| + \sigma \sqrt{2 \log \frac{2n}{\delta}} (2\epsilon + \epsilon_{n,2\delta}) \\ &\leq sB\sigma \sqrt{2 \frac{\log N_{[]} + \log \frac{1}{\delta}}{n}} + \sigma \sqrt{2 \log \frac{2n}{\delta}} (2\epsilon + \epsilon_{n,\delta}) \\ &\leq sB\sigma \sqrt{2 \frac{sK^{**}(2sB)^{1/2} \log \frac{1}{\delta}}{n\epsilon^{1/2}}} + \sigma \sqrt{2 \log \frac{2n}{\delta}} (2\epsilon + \epsilon_{n,\delta}) \\ &\leq sB\sigma \sqrt{2 \frac{sK^{**}(2sB)^{1/2} \log \frac{1}{\delta}}{n\epsilon^{1/2}}} + 2\sigma \sqrt{2 \log \frac{2n}{\delta}} \epsilon + sB\sigma \sqrt{2 \frac{sK^{**}(2sB)^{1/2} \log \frac{1}{\delta}}{n\epsilon^{1/2}}} \log \frac{2n}{\delta} \\ &\leq 2\sigma \sqrt{2 \log \frac{2n}{\delta}} \epsilon + 2sB\sigma \sqrt{\frac{sK^{**}(2sB)^{1/2} \log^2 \frac{2n}{\delta}}{n\epsilon^{1/2}}}. \end{aligned}$$

To balance the two terms, we set $\epsilon = sB \sqrt{\frac{(sK^{**}(sB)^{1/2})^{4/5}}{n^{4/5}}}$. It is easy to verify that if $n \geq c_1 s \sqrt{sB}$ for some absolute constant c_1 , then $\epsilon \in (0, sB\epsilon_3]$ for some absolute constant ϵ_3 as required by the bracketing number results (Corollary 8.3). Furthermore, conditions (8.11) also hold.

In summary, we have that probability at least $1 - \delta$,

$$|\sup_{h \in \mathcal{G}} \langle w, h \rangle| \leq csB\sigma \sqrt{\frac{s^{6/5} B^{2/5} \log^2 \frac{Cn}{\delta}}{n^{4/5}}} \leq csB\sigma \sqrt{\frac{s(sB)^{1/2} \log^2 \frac{Cn}{\delta}}{n^{4/5}}}$$

where we absorbed K^{**} into the constant c and the union bound multipliers into the constant C .

Plugging this result into equation (8.7) we get that, with probability at least $1 - 2\delta$,

$$\begin{aligned}
 \|f_0 - \hat{f}\|_n^2 - \|f_0 - f^*\|_n^2 &\leq csB\sigma\sqrt{\frac{s(sB)^{1/2}\log^2\frac{Cn}{\delta}}{n^{4/5}}} + \lambda 2sB + c(sB)^2\frac{1}{n}\log\frac{2}{\delta} \\
 \|f_0 - \hat{f}\|_n^2 - \|f_0 - f^*\|_n^2 &\leq csB\sigma\sqrt{\frac{s(sB)^{1/2}\log^2\frac{Cn}{\delta}}{n^{4/5}}} + \lambda 2sB \\
 (8.8) \qquad \qquad \qquad &\leq cB\sigma\sqrt{\frac{s^4B^{1/2}}{n^{4/5}}\log^2\frac{Cn}{\delta}} + \lambda 2sB
 \end{aligned}$$

Step 3. Continuing from equation (8.8), we use Lemma 8.1 and another union bound to obtain that, with probability at least $1 - 3\delta$,

$$\begin{aligned}
 \|f_0 - \hat{f}\|_P^2 - \|f_0 - f^*\|_P^2 &\leq cB^2\sigma\sqrt{\frac{s^4}{n^{4/5}}\log^2\frac{Cn}{\delta}} + \lambda 2sB + cB^3\sqrt{\frac{s^5}{n^{4/5}}\log\frac{2}{\delta}} \\
 &\leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log^2\frac{Cn}{\delta}} + \lambda 2sB
 \end{aligned}$$

Substituting in $\lambda \leq 512s\tilde{\sigma}\sqrt{\frac{1}{n}\log^2 np}$ and $\delta = \frac{C}{n}$ we obtain the statement of the theorem. \square

THEOREM 8.4. *Let \hat{g}_k denote the minimizer of the concave postprocessing step with $\lambda_n \leq 512s\tilde{\sigma}\sqrt{\frac{1}{n}\log^2 np}$. Let $\tilde{\sigma} \equiv \max(\sigma, B)$. Suppose n is sufficiently large that $\frac{n^{4/5}}{\log^2 np} \geq c'B^4\tilde{\sigma}^2s^5$ where $c' \geq 1$ is a constant. Then with probability at least $1 - \frac{C}{n}$, for all $k = 1, \dots, s$,*

$$\|f_0 - f^* - \hat{g}_k\|_P^2 - \|f_0 - f^* - g_k^*\|_P^2 \leq cB^2\tilde{\sigma}^{1/2}\sqrt[4]{\frac{s^5}{n^{4/5}}\log^2 np}.$$

PROOF. This proof is similar to that of Theorem 8.3; it requires a few more steps because \hat{g}_k is fitted against $f_0 - \hat{f}$ instead of $f_0 - f^*$. We start with the following decomposition:

$$\begin{aligned}
 \|f_0 - f^* - \hat{g}_k\|_P^2 - \|f_0 - f^* - g_k^*\|_P^2 &= \underbrace{\|f_0 - \hat{f} - \hat{g}_k\|_P^2 - \|f_0 - \hat{f} - g_k^*\|_P^2}_{\text{term 1}} + \\
 &\quad \underbrace{\|f_0 - f^* - \hat{g}_k\|_P^2 - \|f_0 - \hat{f} - \hat{g}_k\|_P^2}_{\text{term 2}} + \\
 (8.9) \qquad \qquad \qquad &\quad \underbrace{\|f_0 - \hat{f} - g_k^*\|_P^2 - \|f_0 - f^* - g_k^*\|_P^2}_{\text{term 3}}.
 \end{aligned}$$

We now bound each of the terms. The proof proceeds almost identically to that of Theorem 8.3, because convex and concave functions have the same bracketing number.

Step 1. To bound term 1, we start from the definition of \hat{g}_k and obtain

$$\begin{aligned} \|y - \hat{f} - \hat{g}_k\|_n^2 + \lambda_n \|\hat{g}\|_\infty &\leq \|y - \hat{f} - g_k^*\|_n^2 + \lambda_n \|g^*\|_\infty \\ \|y - \hat{f} - \hat{g}_k\|_n^2 &\leq \|y - \hat{f} - g_k^*\|_n^2 + \lambda_n 2B \\ \|f_0 - \hat{f} - \hat{g}_k + w\|_n^2 &\leq \|f_0 - \hat{f} - g_k^* + w\|_n^2 + \lambda_n 2B \\ \|f_0 - \hat{f} - \hat{g}_k\|_n^2 - \|f_0 - \hat{f} - g_k^*\|_n^2 &\leq 2\langle w, \hat{g}_k - g_k^* \rangle_n + \lambda_n 2B. \end{aligned}$$

Using the same bracketing analysis as in Step 2 of the proof of Theorem 8.3 but setting $s = 1$, we have, with probability at least $1 - \delta$,

$$\|f_0 - \hat{f} - \hat{g}_k\|_n^2 - \|f_0 - \hat{f} - g_k^*\|_n^2 \leq cB^2\sigma\sqrt{\frac{1}{n^{4/5}}\log\frac{C}{\delta}} + \lambda_n 2B.$$

The condition $n \geq c_1 s \sqrt{sB}$ in the proof of Theorem 8.3 is satisfied here because we assume that $n^{4/5} \geq c_1 B^4 \tilde{\sigma}^2 s^5 \log^2 np$ in the statement of the theorem. Using the uniform convergence result of Lemma 8.1, with probability at least $1 - \delta$,

$$\begin{aligned} \|f_0 - \hat{f} - \hat{g}_k\|_P^2 - \|f_0 - \hat{f} - g_k^*\|_P^2 &\leq cB^2\sigma\sqrt{\frac{1}{n}\log\frac{Cn}{\delta}} + \lambda_n 2B + cB^3\sqrt{\frac{s^5}{n^{4/5}}\log\frac{2}{\delta}} \\ &\leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log\frac{C}{\delta}} + \lambda_n 2B \end{aligned}$$

Finally, plugging in $\lambda_n \leq 9s\tilde{\sigma}\sqrt{\frac{1}{n}\log^2 np}$, we obtain

$$\begin{aligned} \|f_0 - \hat{f} - \hat{g}_k\|_P^2 - \|f_0 - \hat{f} - g_k^*\|_P^2 &\leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log\frac{C}{\delta}} + 2sB\tilde{\sigma}\sqrt{\frac{1}{n}\log^2 np} \\ \|f_0 - \hat{f} - \hat{g}_k\|_P^2 - \|f_0 - \hat{f} - g_k^*\|_P^2 &\leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log^2\frac{Cnp}{\delta}} \end{aligned}$$

with probability at least $1 - \delta$.

Step 2. We now bound term 3.

$$\begin{aligned} \|f_0 - \hat{f} - g_k^*\|_P^2 - \|f_0 - f^* - g_k^*\|_P^2 &\leq \|f_0 - \hat{f}\|_P^2 - \|f_0 - f^*\|_P^2 - 2\langle f_0 - \hat{f}, g_k^* \rangle_P + 2\langle f_0 - f^*, g_k^* \rangle_P \\ &\leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log^2 np} + 2|\langle \hat{f} - f^*, g_k^* \rangle_P| \\ &\leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log^2 np} + 2\|\hat{f} - f^*\|_P \|g_k^*\|_P \\ &\leq cB^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log^2 np} + cB\sqrt{B^2\tilde{\sigma}\sqrt{\frac{s^5}{n^{4/5}}\log^2 np}} \\ &\leq cB^2\tilde{\sigma}^{1/2}\sqrt[4]{\frac{s^5}{n^{4/5}}\log^2 np} \end{aligned}$$

with probability at least $1 - \frac{C}{n}$, by Theorem 8.3. To obtain the fourth inequality, we used the fact that $\|\hat{f} - f^*\|^2 \leq \|f_0 - \hat{f}\|_P^2 - \|f_0 - f^*\|_P^2$, which follows from the fact that f^* is the projection of f_0 onto the set of additive convex functions and the set of additive convex functions is convex itself. The last inequality holds because there is a condition in the theorem which states n is large enough such that $B^2 \tilde{\sigma} \sqrt{\frac{s^5}{n^{4/5}} \log^2 np} \leq 1$. The same derivation and the same bound likewise holds for term 2.

Step 3. Collecting the results and plugging them into equation (8.9), we have, with probability at least $1 - 2\delta$:

$$\|f_0 - f^* - \hat{g}_k\|_P^2 - \|f_0 - f^* - g_k^*\|_P^2 \leq cB^2 \tilde{\sigma}^{1/2} \sqrt[4]{\frac{s^5}{n^{4/5}} \log^2 \frac{4np}{\delta}}$$

Taking a union bound across the s dimensions completes the result. \square

8.3.3. Support Lemmas.

LEMMA 8.1. *Suppose $n \geq c_1 s \sqrt{sB}$ for some absolute constant c_1 . Then, with probability at least $1 - \delta$:*

$$\sup_{f \in \mathcal{C}_B^s} \left| \|f_0 - f\|_n^2 - \|f_0 - f\|_P^2 \right| \leq cB^3 \sqrt{\frac{s^5}{n^{4/5}} \log \frac{2}{\delta}}$$

where c_1 is some absolute constant and c, C are constants possibly dependent on b .

PROOF. Let \mathcal{G} denote the off-centered set of convex functions, that is, $\mathcal{G} \equiv \mathcal{C}^s - f_0$. Note that if $h \in \mathcal{G}$, then $\|h\|_\infty = \|f_0 - f\|_\infty \leq 4sB$. There exists an ϵ -bracketing of \mathcal{G} , and by Corollary 8.3, the bracketing has size at most $\log N_{[]} (2\epsilon, \mathcal{C}^s, L_1(P)) \leq sK^{**} \left(\frac{2sB}{\epsilon} \right)^{1/2}$. By Corollary 8.4, we know that with probability at least $1 - \delta$, $\|h_U - h_L\|_{L_1(P_n)} \leq \epsilon + \epsilon_{n,\delta}$ for all pairs (h_U, h_L) in the bracketing, where $\epsilon_{n,\delta} = sB \sqrt{\frac{K^{**}(2sB)^{1/2} \log \frac{2}{\delta}}{2\epsilon^{1/2}n}}$. Corollary 8.4 necessitates $\epsilon \in (0, sB\epsilon_3]$ for some absolute constant ϵ_3 ; we will verify that this condition holds for large enough n when we set the actual value of ϵ . For a particular function $h \in \mathcal{G}$, we can construct $\psi_L \equiv \min(|h_U|, |h_L|)$ and $\psi_U \equiv \max(|h_U|, |h_L|)$ so that

$$\psi_L^2 \leq h^2 \leq \psi_U^2.$$

We can then bound the $L_1(P)$ norm of $\psi_U^2 - \psi_L^2$ as

$$\begin{aligned} \int (\psi_U^2(x) - \psi_L^2(x))p(x)dx &\leq \int |h_U^2(x) - h_L^2(x)|p(x)dx \\ &\leq \int |h_U(x) - h_L(x)| |h_U(x) + h_L(x)|p(x)dx \\ &\leq 2sB\epsilon \end{aligned}$$

Now we can bound $\|h\|_n^2 - \|h\|_P^2$ as

$$(8.10) \quad \frac{1}{n} \sum_{i=1}^n \psi_L(X_i)^2 - \mathbb{E}\psi_U(X)^2 \leq \|h\|_n^2 - \|h\|_P^2 \leq \frac{1}{n} \sum_{i=1}^n \psi_U(X_i)^2 - \mathbb{E}\psi_L(X)^2$$

Since $\psi_L(X_i)^2$ and $\psi_U(X_i)^2$ are bounded random variables with upper bound $(sB)^2$, Hoeffding's inequality and union bound give that, with probability at least $1 - \delta$, for all ψ_L (and likewise ψ_U)

$$\left| \frac{1}{n} \sum_{i=1}^n \psi_L(X_i)^2 - \mathbb{E}\psi_L(X)^2 \right| \leq (sB)^2 \sqrt{\frac{sK^{**}(sB)^{1/2} \log \frac{2}{\delta}}{\epsilon^{1/2}n}}$$

Plugging this into equation (8.10) above, we have that:

$$\begin{aligned} & \mathbb{E}\psi_L(X)^2 - \mathbb{E}\psi_U(X)^2 - (sB)^2 \sqrt{\frac{sK^{**}(sB)^{1/2} \log \frac{2}{\delta}}{\epsilon^{1/2}n}} \\ & \leq \|h\|_n^2 - \|h\|_P^2 \leq \mathbb{E}\psi_U(X)^2 - \mathbb{E}\psi_L(X)^2 + (sB)^2 \sqrt{\frac{sK^{**}(sB)^{1/2} \log \frac{2}{\delta}}{\epsilon^{1/2}n}}. \end{aligned}$$

Using the $L_1(P)$ norm of $\psi_U^2 - \psi_L^2$ result, we have

$$-sB\epsilon - (sB)^2 \sqrt{\frac{sK^{**}(sB)^{1/2} \log \frac{2}{\delta}}{\epsilon^{1/2}n}} \leq \|h\|_n^2 - \|h\|_P^2 \leq sB\epsilon + (sB)^2 \sqrt{\frac{sK^{**}(sB)^{1/2} \log \frac{2}{\delta}}{\epsilon^{1/2}n}}$$

We balance the terms by choosing $\epsilon = \left(\frac{(sB)^2 sK^{**}(sB)^{1/2}}{n} \right)^{2/5}$. One can easily verify that $\epsilon \leq sB\epsilon_3$ condition needed by Corollary 8.4 is satisfied when $n \geq c_1 s \sqrt{sB}$ for some absolute constant c_1 . We have then that, with probability at least $1 - \delta$,

$$\sup_{h \in \mathcal{G}} \left| \|h\|_n^2 - \|h\|_P^2 \right| \leq cB^3 \sqrt{\frac{s^5 \log \frac{2}{\delta}}{n^{4/5}}}$$

The theorem follows immediately. □

LEMMA 8.2. *Let f_0 and f^* be defined as in Section 8.3.1. Define $\bar{f}^* = \frac{1}{n} \sum_{i=1}^n f^*(X_i)$. Then, with probability at least $1 - 2\delta$,*

$$\left| \|f_0 - f^*\|_n^2 - \|f_0 - f^* + \bar{f}^*\|_n^2 \right| \leq c(sB)^2 \frac{1}{n} \log \frac{4}{\delta}$$

PROOF. We decompose the empirical norm as

$$\begin{aligned} \|f_0 - f^* + \bar{f}^*\|_n^2 &= \|f_0 - f^*\|_n^2 + 2\langle f_0 - f^*, \bar{f}^* \rangle + \bar{f}^{*2} \\ &= \|f_0 - f^*\|_n^2 + 2\bar{f}^* \langle f_0 - f^*, \mathbf{1} \rangle_n + \bar{f}^{*2} \\ &= \|f_0 - f^*\|_n^2 + 2\bar{f}^* \bar{f}_0 - \bar{f}^{*2}. \end{aligned}$$

Now $\bar{f}^* = \frac{1}{n} \sum_{i=1}^n f^*(X_i)$ is the average of n bounded mean-zero random variables and therefore, with probability at least $1 - \delta$, $|\bar{f}^*| \leq 4sB\sqrt{\frac{1}{n} \log \frac{2}{\delta}}$. The same reasoning likewise applies to $\bar{f}_0 = \frac{1}{n} \sum_{i=1}^n f_0(X_i)$.

Taking a union bound and we have that, with probability at least $1 - 2\delta$,

$$\begin{aligned} |\bar{f}^*| |\bar{f}_0| &\leq c(sB)^2 \frac{1}{n} \log \frac{2}{\delta} \\ \bar{f}^{*2} &\leq c(sB)^2 \frac{1}{n} \log \frac{2}{\delta} \end{aligned}$$

Therefore, with probability at least $1 - 2\delta$,

$$\|f_0 - f^*\|_n^2 - c(sB)^2 \frac{1}{n} \log \frac{2}{\delta} \leq \|f_0 - f^* + \bar{f}^*\|_n^2 \leq \|f_0 - f^*\|_n^2 + c(sB)^2 \frac{1}{n} \log \frac{2}{\delta}$$

□

8.3.4. Supporting Lemma for Theorem 3.2.

LEMMA 8.3. *Let $f : [0, 1]^p \rightarrow \mathbb{R}$ be a twice differentiable function. Suppose $p(\mathbf{x})$ is a density on $[0, 1]^p$ such that $\partial_{x_k} p(\mathbf{x}_{-k} | x_k)$ and $\partial_{x_k}^2 p(\mathbf{x}_{-k} | x_k)$ are continuous as functions of x_k . Let $\phi(\mathbf{x}_{-k})$ be a continuous function not dependent on x_k .*

Then, $h_k^(x_k) \equiv \mathbb{E}[f(X) - \phi(X_{-k}) | x_k]$ is twice differentiable and has a second derivative lower bounded away from $-\infty$.*

PROOF. We can write

$$h_k^*(x_k) = \int_{\mathbf{x}_{-k}} (f(\mathbf{x}) - \phi(\mathbf{x}_{-k})) p(\mathbf{x}_{-k} | x_k) d\mathbf{x}_{-k}$$

The integrand is bounded because it is a sum-product of continuous functions over a compact set. Therefore, we can differentiate under the integral and derive that

$$\begin{aligned} \partial_{x_k} h_k^*(x_k) &= \int_{\mathbf{x}_{-k}} f'(\mathbf{x}) p(\mathbf{x}_{-k} | x_k) + (f(\mathbf{x}) - \phi(\mathbf{x}_{-k})) p'(\mathbf{x}_{-k} | x_k) d\mathbf{x}_{-k} \\ \partial_{x_k}^2 h_k^*(x_k) &= \int_{\mathbf{x}_{-k}} f''(\mathbf{x}) p(\mathbf{x}_{-k} | x_k) + 2f'(\mathbf{x}) p'(\mathbf{x}_{-k} | x_k) + (f(\mathbf{x}) - \phi(\mathbf{x}_{-k})) p''(\mathbf{x}_{-k} | x_k) d\mathbf{x}_{-k} \end{aligned}$$

where we have used the shorthand $f'(\mathbf{x})$, $p'(\mathbf{x}_{-k} | x_k)$ to denote $\partial_{x_k} f(\mathbf{x})$, $\partial_{x_k} p(\mathbf{x}_{-k} | x_k)$, etc.

This proves that $h_k^*(x_k)$ is twice-differentiable. To see that the second derivative is lower bounded, we note that $f''(\mathbf{x}) p(\mathbf{x}_{-k} | x_k)$ is non-negative and the other terms in the second-derivative are all continuous functions on a compact set and thus bounded. □

LEMMA 8.4. *Let $p(\mathbf{x})$ be a positive density over $[0, 1]^p$. Let $\phi(\mathbf{x}) = \sum_{j=1}^p \phi_j(x_j)$ be an additive function.*

Suppose that for all j , $\phi_j(x_j)$ is bounded and $\mathbb{E}\phi_j(X_j) = 0$. Suppose for some j , $\mathbb{E}\phi_j(X_j)^2 > 0$, then it must be that $\mathbb{E}\phi(X)^2 > 0$.

PROOF. Suppose, for sake of contradiction, that

$$\mathbb{E}\phi(X)^2 = \mathbb{E}(\phi_j(X_j) + \phi_{-j}(X_{-j}))^2 = 0.$$

Let $A_+ = \{x_j : \phi_j(x_j) \geq 0\}$. Since $\mathbb{E}\phi_j(X_j) = 0$, $\mathbb{E}\phi_j(X_j)^2 > 0$, and ϕ_j is bounded, it must be that both A_+ has probability greater than 0.

Now, define $B_+ = \{x_{-j} : \phi_{-j}(x_{-j}) \geq 0\}$. B_+ then must have probability greater than 0 as well.

Since $p(\mathbf{x})$ is a positive density, the set $A_+ \times B_+$ must have a positive probability. However, $\phi > 0$ on $A_+ \times B_+ \subset [0, 1]^p$ which implies $\mathbb{E}\phi(X)^2 > 0$. □

8.3.5. *Concentration of Measure.* A sub-exponential random is the square of a sub-Gaussian random variable [Vershynin \(2010\)](#).

PROPOSITION 8.1. (*Subexponential Concentration* [Vershynin \(2010\)](#)) Let X_1, \dots, X_n be zero-mean independent subexponential random variables with subexponential scale K . Then

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i\right| \geq \epsilon\right) \leq 2 \exp\left[-cn \min\left(\frac{\epsilon^2}{K^2}, \frac{\epsilon}{K}\right)\right]$$

where $c > 0$ is an absolute constant.

For uncentered subexponential random variables, we can use the following fact. If X_i subexponential with scale K , then $X_i - \mathbb{E}[X_i]$ is also subexponential with scale at most $2K$. Restating, we can set

$$c \min\left(\frac{\epsilon^2}{K^2}, \frac{\epsilon}{K}\right) = \frac{1}{n} \log \frac{1}{\delta}.$$

Thus, with probability at least $1 - \delta$, the deviation is at most

$$K \max\left(\sqrt{\frac{1}{cn} \log \frac{C}{\delta}}, \frac{1}{cn} \log \frac{C}{\delta}\right).$$

COROLLARY 8.1. Let W_1, \dots, W_n be n independent sub-Gaussian random variables with sub-Gaussian scale σ . Then, for all $n > n_0$, with probability at least $1 - \frac{1}{n}$,

$$\frac{1}{n} \sum_{i=1}^n W_i^2 \leq c\sigma^2.$$

PROOF. Using the subexponential concentration inequality, we know that, with probability at least $1 - \frac{1}{n}$,

$$\left|\frac{1}{n} \sum_{i=1}^n W_i^2 - \mathbb{E}W^2\right| \leq \sigma^2 \max\left(\sqrt{\frac{1}{cn} \log \frac{C}{\delta}}, \frac{1}{cn} \log \frac{C}{\delta}\right).$$

First, let $\delta = \frac{1}{n}$. Suppose n is large enough such that $\frac{1}{cn} \log Cn < 1$. Then, we have, with probability at least $1 - \frac{1}{n}$,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n W_i^2 &\leq c\sigma^2 \left(1 + \sqrt{\frac{1}{cn} \log Cn}\right) \\ &\leq 2c\sigma^2. \end{aligned}$$

□

8.3.6. Sampling Without Replacement.

LEMMA 8.5. ([Serfling \(1974\)](#)) Let x_1, \dots, x_N be a finite list, $\bar{x} = \mu$. Let X_1, \dots, X_n be sampled from x without replacement.

Let $b = \max_i x_i$ and $a = \min_i x_i$. Let $r_n = 1 - \frac{n-1}{N}$. Let $S_n = \sum_i X_i$. Then we have that

$$\mathbb{P}(S_n - n\mu \geq n\epsilon) \leq \exp\left(-2n\epsilon^2 \frac{1}{r_n(b-a)^2}\right).$$

COROLLARY 8.2. Suppose $\mu = 0$.

$$\mathbb{P}\left(\frac{1}{N}S_n \geq \epsilon\right) \leq \exp\left(-2N\epsilon^2 \frac{1}{(b-a)^2}\right)$$

And, by union bound, we have that

$$\mathbb{P}\left(\left|\frac{1}{N}S_n\right| \geq \epsilon\right) \leq 2 \exp\left(-2N\epsilon^2 \frac{1}{(b-a)^2}\right)$$

A simple restatement is that with probability at least $1 - \delta$, the deviation $|\frac{1}{N}S_n|$ is at most $(b-a)\sqrt{\frac{1}{2N} \log \frac{2}{\delta}}$.

PROOF.

$$\mathbb{P}\left(\frac{1}{N}S_n \geq \epsilon\right) = \mathbb{P}\left(S_n \geq \frac{N}{n}n\epsilon\right) \leq \exp\left(-2n \frac{N^2}{n^2} \epsilon^2 \frac{1}{r_n(b-a)^2}\right).$$

We note that $r_n \leq 1$ always, and $n \leq N$ always. Thus,

$$\exp\left(-2n \frac{N^2}{n^2} \epsilon^2 \frac{1}{r_n(b-a)^2}\right) \leq \exp\left(-2N\epsilon^2 \frac{1}{(b-a)^2}\right)$$

completing the proof. □

8.3.7. *Bracketing Numbers for Convex Functions.*

DEFINITION 8.1. Let \mathcal{C} be a set of functions. For a given ϵ and metric ρ (which we take to be L_2 or $L_2(P)$), we define a *bracketing* of \mathcal{C} to be a set of pairs of functions $\{(f_L, f_U)\}$ satisfying (1) $\rho(f_L, f_U) \leq \epsilon$ and (2) for any $f \in \mathcal{C}$, there exist a pair (f_L, f_U) where $f^U \geq f \geq f^L$.

We let $N_{[]}(\epsilon, \mathcal{C}, \rho)$ denote the size of the smallest bracketing of \mathcal{C}

PROPOSITION 8.2. (*Proposition 16 in Kim and Samworth (2014)*) Let \mathcal{C} be the set of convex functions supported on $[-1, 1]^d$ and uniformly bounded by B . Then there exist constants ϵ_3 and K^{**} , dependent on d , such that

$$\log N_{[]} (2\epsilon, \mathcal{C}, L_2) \leq K^{**} \left(\frac{2B}{\epsilon} \right)^{d/2}$$

for all $\epsilon \in (0, B\epsilon_3]$.

It is trivial to extend Kim and Samworth's result to the $L_2(P)$ norm for an absolutely continuous distribution P .

PROPOSITION 8.3. Let P be a distribution with a density p . Let $\mathcal{C}, B, \epsilon_3, K^{**}$ be defined as in Proposition 8.2. Then,

$$\log N_{[]} (2\epsilon, \mathcal{C}, L_1(P)) \leq K^{**} \left(\frac{2B}{\epsilon} \right)^{d/2}$$

for all $\epsilon \in (0, B\epsilon_3]$.

PROOF. Let \mathcal{C}_ϵ be the bracketing that satisfies the size bound in Proposition 8.3. Let $(f_L, f_U) \in \mathcal{C}_\epsilon$. Then we have that:

$$\begin{aligned} \|f_L - f_U\|_{L_1(P)} &= \int |f_L(x) - f_U(x)| p(x) dx \\ &\leq \left(\int |f_L(x) - f_U(x)|^2 dx \right)^{1/2} \left(\int p(x)^2 dx \right)^{1/2} \\ &\leq \left(\int |f_L(x) - f_U(x)|^2 dx \right)^{1/2} \\ &\leq \|f_L - f_U\|_{L_2} \leq \epsilon. \end{aligned}$$

On the third line, we used the fact that $\int p(x)^2 dx \leq (\int p(x) dx)^2 \leq 1$. \square

It is also simple to extend the bracketing number result to additive convex functions. As before, let \mathcal{C}^s be the set of additive convex functions with s components.

COROLLARY 8.3. *Let P be a distribution with a density p . Let B, ϵ_3, K^{**} be defined as in Proposition 8.2. Then,*

$$\log N_{[]} (2\epsilon, \mathcal{C}^s, L_1(P)) \leq sK^{**} \left(\frac{2sB}{\epsilon} \right)^{1/2}$$

for all $\epsilon \in (0, sB\epsilon_3]$.

PROOF. Let $f \in \mathcal{C}^s$. We can construct an ϵ -bracketing for f through ϵ/s -bracketings for each of the components $\{f_k\}_{k=1, \dots, s}$:

$$f_U = \sum_{k=1}^s f_{Uk} \quad f_L = \sum_{k=1}^s f_{Lk}$$

It is clear that $f_U \geq f \geq f_L$. It is also clear that $\|f_U - f_L\|_{L_1(P)} \leq \sum_{k=1}^s \|f_{Uk} - f_{Lk}\|_{L_1(P)} \leq \epsilon$. \square

The following result follows from Corollary 8.3 directly by a union bound.

COROLLARY 8.4. *Let X_1, \dots, X_n be random samples from a distribution P . Let $1 > \delta > 0$. Let \mathcal{C}_ϵ^s be an ϵ -bracketing of \mathcal{C}^s with respect to the $L_1(P)$ -norm whose size is at most $N_{[]} (2\epsilon, \mathcal{C}^s, L_1(P))$. Let $\epsilon \in (0, sB\epsilon_3]$.*

Then, with probability at least $1 - \delta$, for all pairs $(f_L, f_U) \in \mathcal{C}_\epsilon^s$, we have that

$$\frac{1}{n} \sum_{i=1}^n |f_L(X_i) - f_U(X_i)| \leq \epsilon + \epsilon_{n,\delta}$$

where

$$\epsilon_{n,\delta} \equiv sB \sqrt{\frac{\log N_{[]} (2\epsilon, \mathcal{C}^s, L_2(P)) + \log \frac{1}{\delta}}{2n}} = \sqrt{\frac{sK^{**}(sB)^{1/2}}{2\epsilon^{1/2}n} + \frac{1}{2n} \log \frac{1}{\delta}}.$$

PROOF. Noting that $|f_L(X_i) - f_U(X_i)|$ is at most sB and there are $N_{[]} (2\epsilon, \mathcal{C}^s, L_1(P))$ pairs (f_L, f_U) , the inequality follows from a direct application of a union bound and Hoeffding's Inequality. \square

To make the expression in this corollary easier to work with, we derive an upper bound for $\epsilon_{n,\delta}$. Suppose

$$(8.11) \quad \epsilon^{1/2} \leq 2sK^{**}(sB)^{1/2} \quad \text{and} \quad \log \frac{1}{\delta} \geq 2.$$

Then we have that

$$\epsilon_n \leq sB \sqrt{\frac{sK^{**}(sB)^{1/2} \log \frac{1}{\delta}}{\epsilon^{1/2}n}}.$$

9. Gaussian Example. Let H be a positive definite matrix and let $f(x_1, x_2) = H_{11}x_1^2 + 2H_{12}x_1x_2 + H_{22}x_2^2 + c$ be a quadratic form where c is a constant such that $\mathbb{E}[f(X)] = 0$. Let $X \sim N(0, \Sigma)$ be a random bivariate Gaussian vector with covariance $\Sigma = [1, \alpha; \alpha, 1]$

PROPOSITION 9.1. *Let $f_1^*(x_1) + f_2^*(x_2)$ be the additive projection of f under the bivariate Gaussian distribution. That is,*

$$f_1^*, f_2^* \equiv \arg \min_{f_1, f_2} \left\{ \mathbb{E} (f(X) - f_1(X_1) - f_2(X_2))^2 : \mathbb{E}[f_1(X_1)] = \mathbb{E}[f_2(X_2)] = 0 \right\}$$

Then, we have that

$$\begin{aligned} f_1^*(x_1) &= \left(\frac{T_1 - T_2\alpha^2}{1 - \alpha^4} \right) x_1^2 + c_1 \\ f_2^*(x_2) &= \left(\frac{T_2 - T_1\alpha^2}{1 - \alpha^4} \right) x_2^2 + c_2 \end{aligned}$$

where $T_1 = H_{11} + 2H_{12}\alpha + H_{22}\alpha^2$ and $T_2 = H_{22} + 2H_{12}\alpha + H_{11}\alpha^2$ and c_1, c_2 are constants such that $\mathbb{E}[f_1^*(X_1)] = \mathbb{E}[f_2^*(X_2)] = 0$.

PROOF. By Lemma 3.1, we need only verify that f_1^*, f_2^* satisfy

$$\begin{aligned} f_1^*(x_1) &= \mathbb{E}[f(X) - f_2^*(X_2) | x_1] \\ f_2^*(x_2) &= \mathbb{E}[f(X) - f_1^*(X_1) | x_2]. \end{aligned}$$

Let us guess that f_1^*, f_2^* are quadratic forms $f_1^*(x_1) = a_1x_1^2 + c_1$, $f_2^*(x_2) = a_2x_2^2 + c_2$ and verify that there exist a_1, a_2 to satisfy the above equations. Since we are not interested in constants, we define \simeq to be equality up to a constant. Then,

$$\begin{aligned} &\mathbb{E}[f(X) - f_2^*(X_2) | x_1] \\ &\simeq \mathbb{E}[H_{11}X_1^2 + 2H_{12}X_1X_2 + H_{22}X_2^2 - a_2X_2^2 | x_1] \\ &\simeq H_{11}x_1^2 + 2H_{12}x_1\mathbb{E}[X_2 | x_1] + H_{22}\mathbb{E}[X_2^2 | x_1] - a_2\mathbb{E}[X_2^2 | x_1] \\ &\simeq H_{11}x_1^2 + 2H_{12}\alpha x_1^2 + H_{22}\alpha^2 x_1^2 - a_2\alpha^2 x_1^2 \\ &\simeq (H_{11} + 2H_{12}\alpha + H_{22}\alpha^2 - a_2\alpha^2)x_1^2. \end{aligned}$$

Likewise, we have that

$$\mathbb{E}[f(X) - f_1^*(X_1) | x_2] \simeq (H_{22} + 2H_{12}\alpha + H_{11}\alpha^2 - a_1\alpha^2)x_2^2.$$

Thus, a_1, a_2 need only satisfy the linear system

$$\begin{aligned} T_1 - a_2\alpha^2 &= a_1 \\ T_2 - a_1\alpha^2 &= a_2 \end{aligned}$$

where $T_1 = H_{11} + 2H_{12}\alpha + H_{22}\alpha^2$ and $T_2 = H_{22} + 2H_{12}\alpha + H_{11}\alpha^2$. It is then simple to solve the system and verify that a_1, a_2 are as specified. \square