

1 3/1/2014

This note analyzes SCAM where we penalize $\|f\|_\infty$ instead of $\|\partial f\|_\infty$.

1.1 Changes to Optimization

The optimization program must change. We modify two formulations of the optimization. The first is the formulation in which we use both f and β as variables.

$$\begin{aligned} \min_{h_k, \beta_k, \gamma_k} & \frac{1}{2n} \sum_{i=1}^n \left(y_i - \sum_{k' \neq k} h_{k'i} - h_{ki} \right)^2 + \lambda \sum_k \|h_k\|_\infty \\ \text{s.t.} & h_{k(i+1)} = h_{k(i)} + \beta_{k(i)}(x_{k(i+1)} - x_{k(i)}) \\ & \beta_{k(i+1)} \geq \beta_{k(i)}, \mathbf{1}h_k = 0 \end{aligned}$$

Of course, the $\lambda\|h_k\|_\infty$ penalty can be replaced by $\lambda\gamma_k$ and linear inequalities involving γ_k .

Now we reformulate the optimization program that is in term of the discretized second derivative d_k .

$$\begin{aligned} \min_{d_k} & \frac{1}{2n} \|Y - \sum_{k=1}^p \Delta_k d_k\|_2^2 + \lambda \sum_{k=1}^p \|\Delta_k d_k\|_\infty \\ \text{s.t.} & d_{k(2)}, \dots, d_{k(n-1)} \geq 0 \\ & \mathbf{1}_n^\top \Delta_k d_k = 0 \quad \forall k \end{aligned}$$

1.2 Subgradient analysis

We take the subgradient of the optimization program with d_k .

First, we note that the subgradient of the sup-norm $\partial\|x\|_\infty$ at $x = 0$ is $\{z : \sum_i z_i = 1\}$.

We fix one dimension k , let $Y' = Y - \sum_{k' \neq k} \Delta_{k'} d_{k'}$.

The Lagrangian form of the optimization is

$$\mathcal{L}(d_k, v, \gamma) = \frac{1}{2n} \|Y' - \Delta_k d_k\|_2^2 + \lambda \|\Delta_k d_k\|_\infty - \sum_{i=2}^{n-1} v_i d_{ki} + \gamma \mathbf{1}_n^\top \Delta_k d_k$$

with the constraint that $v_i \geq 0$ for all i .

We want to find under what conditions does the solution $d_k = 0$ satisfy the KKT conditions.

$$\partial_{d_k} \mathcal{L} = \frac{1}{n} \Delta_k^\top (Y' - \Delta_k d_k) + \lambda \Delta_k^\top z - v + \gamma \mathbf{1}_n^\top \Delta_k$$

where $z \in \partial\|\Delta_k d_k\|_\infty$.

If we evaluate the subgradient at $d_k = 0$, then we have that

$$\partial_{d_k} \mathcal{L} \Big|_{d_k=0} = \frac{1}{n} \Delta_k^\top Y' + \lambda \Delta_k^\top z - v + \gamma \mathbf{1}_n^\top \Delta_k$$

where $\|z\|_1 \leq 1$.

We want to argue that $d_k = 0$ is an optimal solution. If $d_k = 0$, then complementary slackness and primal feasibility are obvious satisfied. We need only verify stationarity and dual feasibility then.

Let us take a brief digression and understand the Δ_k matrix a little bit more. Δ_k is $n \times n - 1$. Each row corresponds to sample i ; each column corresponds to an order (j) . Each entry (i, j) is $[X_{ki} - X_{k(j)}]_+$. Let us reorder the samples so that the i -th sample is the i -smallest sample.

We will construct $\gamma = 0$, and $z = (0, 0, \dots, a)$ for some $0 < a < 1$. (coordinates of z correspond to the new sample ordering) We then just need to show that

$$\begin{aligned} & \frac{1}{n} \Delta_k^\top Y' + \lambda \Delta_k^\top z \geq 0 \\ & \frac{1}{n} \sum_{i>j} (X_{ki} - X_{kj}) Y'_i + \lambda (X_{kn} - X_{kj}) a \geq 0 \quad \text{for each } j \\ & \frac{1}{n} \sum_{i>j} \sum_{j<i' \leq i} \text{gap}_{i'} Y'_i + \lambda (X_{kn} - X_{kj}) a \geq 0 \\ & \frac{1}{n} \sum_{i'>j} \text{gap}_{i'} \sum_{i \geq i'} Y'_i + \lambda (X_{kn} - X_{kj}) a \geq 0 \\ & \frac{1}{n} \sum_{i'>j} \text{gap}_{i'} \mathbf{1}_{i':n}^\top Y' + \lambda (X_{kn} - X_{kj}) a \geq 0 \end{aligned}$$

Where $\text{gap}_i = X_{ki} - X_{k,i-1}$. (with respect to ordered indices) The notation $\mathbf{1}_{i':n}$ is a vector that is one on the last i' to n coordinates and zero elsewhere.

Suppose we show that $\frac{1}{n} \mathbf{1}_{i':n}^\top Y'$ is on the order of $O(\frac{1}{\sqrt{n}})$, then we would be done. Here, we will need to bound $\|Y'\|_\infty$.

$Y'_i = Y_i - \hat{f}(X_i) = f_0(X_i) + w_i - \hat{f}(X_i)$. Both $f_0(X_i)$ and $\hat{f}(X_i)$ are bounded by $2sLb$.

For a single w_i , we know that with probability at least $1 - \delta$, $|w_i| \leq c\sigma \sqrt{\log \frac{2}{\delta}}$. Therefore, taking an union bound and setting $\delta = 1/n$, we have that $\max_i |w_i| \leq c\sigma \sqrt{\log n}$ with probability at least $1 - \frac{1}{n}$.

Putting this together, $\|Y'\|_\infty \leq c(sLb + \sigma \sqrt{\log n}) \leq csLb\sigma \sqrt{\log n}$.

Now, we also need to bound the mean of Y' :

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n f_0(X_i) + w_i - \hat{f}(X_i) \\ & = \frac{1}{n} \sum_{i=1}^n f_0(X_i) + w_i \end{aligned}$$

The first mean $\frac{1}{n} \sum_{i=1}^n f_0(X_i)$ is the average of n bounded zero-mean random variables since $|f_0(X_i)| \leq 2sLb$. Therefore, the first term is at most $sLb \sqrt{\frac{1}{n} \log \frac{2}{\delta}}$ with probability at least $1 - \delta$.

The second mean is the average of n subgaussian random variables with subgaussian scale σ . Therefore, the second term is at most $\sigma \sqrt{\frac{1}{n} \log \frac{2}{\delta}}$.

The mean of Y' is therefore at most $\mathbf{c}(sLb + \sigma)\sqrt{\frac{1}{n} \log n}$ with probability at least $1 - \frac{1}{n}$.
By Serfling and union bound, with probability at least $1 - \frac{1}{n}$, we have that the deviation for all i

$$\begin{aligned} \left| \frac{1}{n} \mathbf{1}_{i:n}^\top Y' \right| &\leq \mathbf{c}sLb\sigma\sqrt{\log n}\sqrt{\frac{1}{n} \log 2np} + \mathbf{c}(sLb + \sigma)\sqrt{\frac{1}{n} \log np} \\ &\leq \mathbf{c}sLb\sigma\sqrt{\frac{1}{n} \log n \log np} \end{aligned}$$

Plugging this in and plugging in $\lambda > \mathbf{c}sLb\sigma\sqrt{\frac{1}{n} \log n \log np}$, we have that:

$$\begin{aligned} \frac{1}{n} \sum_{i' > j} \mathbf{gap}_{i'} \mathbf{1}_{i':n}^\top Y' + \lambda(X_{kn} - X_{kj})a &\geq \lambda(X_{kn} - X_{kj})a - \sum_{i' > j} \mathbf{gap}_{i'} \mathbf{c}sLb\sigma\sqrt{\frac{1}{n} \log n \log np} \\ &\geq \lambda(X_{kn} - X_{kj})a - (X_{kn} - X_{kj})\mathbf{c}sLb\sigma\sqrt{\frac{1}{n} \log n \log np} \\ &\geq 0 \end{aligned}$$

1.3 False Negative Analysis

Before we begin, we keep in mind that the analysis should be flexible and should be easy to modify to accomodate the following:

1. choice of norm in the penalty
2. with or without a Lipschitz constraint, boundedness constraint

1.3.1 Notation

Where $f : \mathbb{R}^p \rightarrow \mathbb{R}$, we let $\|f\|_P \equiv \mathbb{E}f(X)^2$.

1.3.2 Preliminary

We start with the definitions. Let $y = (y_1, \dots, y_n)$ and $y_i = f_0(x_i) + w_i$. We assume f_0 to be convex.¹ We assume that all $x_i \in [-b, b]^s$

- Let \mathcal{C}^1 denote the set of univariate convex functions supported on $[-b, b]$. Let $\mathcal{C}_L^1 \equiv \{f \in \mathcal{C}^1 : \|\partial f\|_\infty \leq L\}$ denote the set of L -Lipschitz univariate convex functions.
- Define \mathcal{C}^s as the set of convex additive functions

$$\mathcal{C}^s \equiv \{f : f = \sum_{k=1}^s f_k, f_k \in \mathcal{C}^1\}$$

- We will also define classes of bounded and Lipschitz convex additive functions.

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

- 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 B be an upper bound on $\|f_0\|_\infty$ and $\|f^*\|_\infty$ and let L be an upper bound on $\|\partial_{x_k} f_0\|_\infty$ and $\|\partial f_k^*\|_\infty$.
- We define \hat{f} as the empirical risk minimizer:

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

¹We can probably relax this assumption to be a little bit more general.

- We also define \hat{f}^* as the noiseless empirical risk minimizer:

$$\hat{f}^* = \arg \min \{ \|f_0 - f\|_n^2 : f \in \mathcal{C}_L^s, \mathbf{1}_n^\top f_k = 0 \}$$

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

$$g_k^* \equiv \min_{g_k \in \mathcal{C}} \|f_0 - f_{-k}^* - g_k\|_P^2$$

In our proof, we will analyze g_k^* for k 's such that $f_k^* = 0$.

- Also, we define the noiseless empirical version:

$$\hat{g}_k^* \equiv \min_{g_k \in \mathcal{C}_L} \|f_0 - f_{-k}^* - g_k\|_n^2$$

- Likewise, we define the empirical version

$$\hat{g}_k \equiv \min_{g_k \in \mathcal{C}_L} \|y - \hat{f} - g_k\|_n^2$$

By the definition of the ACDC procedure, \hat{g}_k exist only for k that have zero in their convex additive approximation.

1.3.3 Proof outline

The **final step** is to show the following:

$$\begin{aligned} \|f_0 - \hat{f}\|_P &\approx \|f_0 - f^*\|_P \\ \|f_0 - f^* - \hat{g}_k\|_P &\approx \|f_0 - f^* - g_k^*\|_P \quad \text{for all } k \in \{1, \dots, s\} \text{ where } f_k^* = 0 \end{aligned}$$

where the difference is a term that decreases with n .

Suppose $\hat{f}_k = 0$ and $f_k^* \neq 0$, then for large enough n , there must be a contradiction because $\|f_0 - f^*\|_P > \|f_0 - f^{*(s-1)}\|_P$. We would have a contradiction.

Suppose $f_k^* = 0$, then $g_k^* \neq 0$. We would have another contradiction.

Concavity, **removing** sample-dependent convex functions.

For arbitrary function g , we have the following:

$$\begin{aligned} \|f_0 - \hat{f} - g\|_n &= \|f_0 - \hat{f} + f^* - f^* - g\|_n \leq \|f_0 - f^* - g\|_n + \|\hat{f} - f^*\|_n \\ &\geq \|f_0 - f^* - g\|_n - \|\hat{f} - f^*\|_n \end{aligned}$$

Therefore:

$$\|f_0 - f^* - g\|_n^2 - 2\|f_0 - f^* - g\|_n \|\hat{f} - f^*\|_n \leq \|f_0 - \hat{f} - g\|_n^2 \leq \|f_0 - f^* - g\|_n^2 + 2\|f_0 - f^* - g\|_n \|\hat{f} - f^*\|_n$$

The radius of approximation is $2\|f_0 - f^* - g\|_n \|\hat{f} - f^*\|_n$, which we need to upper bound. $\|f_0 - f^* - g\|_n \leq csLb$.

We **start** from the definition. We don't need to force y to be zero-mean.

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

The $\|f_0 - \hat{f}\|_n^2 - \|f_0 - \hat{f}^*\|_n^2$ term is always positive. We want to lower bound the second term with a negative quantity.

One way to proceed: $\|\hat{f}_k^*\|_\infty \leq 2bL$. And because \hat{f}^* minimizes

Note: The \hat{f}^* here could be any function that does depend on the noise. The advantage of having \hat{f}^* is that $\|f_0 - \hat{f}\|_n^2 - \|f_0 - \hat{f}^*\|_n^2 \geq \|\hat{f} - \hat{f}^*\|_n^2$.

Continuation from start. Taking the **start** derivation and instead of comparing against \hat{f}^* , we will instead compare against f^* (de-measured).

Since the samples X_1, \dots, X_n are clear from context, we denote \bar{f}^* as $\frac{1}{n} \sum_{i=1}^n f^*(X_i)$.

All parts still work and we get down to the following inequality:

$$\begin{aligned} \|f_0 - \hat{f}\|_n^2 - \|f_0 - f^* + \bar{f}^*\|_n^2 &\leq 2|\langle w, \hat{f} - f^* + \bar{f}^* \rangle| + \lambda 2sLb \\ &\leq \mathbf{c}(sLb)^2 \sigma \sqrt{\frac{1}{n} \log 4n} + \lambda 2sLb \end{aligned}$$

with probability at least $1 - \frac{1}{n}$.

We use the *centering effect* analysis. Since with high probability $(1 - \frac{1}{n})$

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

Therefore, with probability at least $1 - \frac{2}{n}$,

$$\begin{aligned} \|f_0 - \hat{f}\|_n^2 - \|f_0 - f^*\|_n^2 &\leq \mathbf{c}(sLb)^2 \frac{1}{n} \log 4n + \mathbf{c}(sLb)^2 \sigma \sqrt{\frac{1}{n} \log 4n} + \lambda 2sLb \\ &\leq \mathbf{c}(Lb)^2 \sigma \sqrt{\frac{s^4}{n} \log 4n} + \lambda 2sLb \end{aligned}$$

Now, we use the the *uniform convergence* result. With probability at least $1 - \frac{3}{n}$

$$\begin{aligned} \|f_0 - \hat{f}\|_P^2 - \|f_0 - f^*\|_P^2 &\leq \mathbf{c}(Lb)^2 \sigma \sqrt{\frac{s^4}{n} \log 4n} + \lambda 2sLb + \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 2n} \\ &\leq \mathbf{c}(Lb)^3 \sigma \sqrt{\frac{s^5}{n^{4/5}} \log 4n} + \lambda 2sLb \end{aligned}$$

f convergence.

Here, our goal is to bound $\|\hat{f} - f^*\|_n^2$.

Since $f^* = \arg \min_{f \in \mathcal{C}^s} \|f_0 - f\|_P^2$, we have that

$$\begin{aligned} \|f^* - \hat{f}\|_P^2 &\leq \|f_0 - \hat{f}\|_P^2 - \|f_0 - f^*\|_P^2 \\ &\leq \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 4n} + \lambda 2sLb \end{aligned}$$

Now, by *uniform convergence* again, we have:

$$\|f^* - \hat{f}\|_n^2 \leq \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 4n} + \mathbf{c}\lambda sLb$$

Centering Effect

Here, we bound the difference of $\|f_0 - f^* + \bar{f}^*\|_n^2$ and $\|f_0 - f^*\|_n^2$.

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

$\bar{f}^* = \frac{1}{n} \sum_{i=1}^n f^*(X_i)$ is the average of n bounded mean-zero random variables and therefore, with probability at least $1 - \delta$, $|\bar{f}^*| \leq 4sLb\sqrt{\frac{1}{n} \log \frac{2}{\delta}}$.

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 letting $\delta = \frac{1}{n}$, we have that, with probability at least $1 - \frac{1}{n}$,

$$\begin{aligned}|\bar{f}^*| |\bar{f}_0| &\leq \mathbf{c}(sLb)^2 \frac{1}{n} \log 4n \\ \bar{f}^{*2} &\leq \mathbf{c}(sLb)^2 \frac{1}{n} \log 4n\end{aligned}$$

Therefore, with probability $1 - \frac{1}{n}$,

$$\|f_0 - f^*\|_n^2 - \mathbf{c}(sLb)^2 \frac{1}{n} \log 4n \leq \|f_0 - f^* + \bar{f}^*\|_n^2 \leq \|f_0 - f^*\|_n^2 + \mathbf{c}(sLb)^2 \frac{1}{n} \log 4n$$

Denoising

Here we outline how to bound $\langle w, \hat{f} - \hat{f}^* \rangle$ where \hat{f}^* can be any function that doesn't depend on w .

We suppose that both $\hat{f}, \hat{f}^* \in \mathcal{C}_L^s$.

By Bronshtein, we know that:

$$\begin{aligned}\log N_\infty(\epsilon, \mathcal{C}_L^1) &= CbL\epsilon^{-1/2} \\ \log N_\infty(\epsilon, \mathcal{C}_L^s) &= CsbLs^{1/2}\epsilon^{-1/2}\end{aligned}$$

Clearly, the covering number would not change if we center the set \mathcal{C}_L^s around \hat{f}^* , call this \mathcal{G} . The difference $\hat{f} - \hat{f}^* \in \mathcal{G}$

For all $h \in \mathcal{G}$:

$$\begin{aligned}|\langle w, h \rangle_n| &= |\langle w, h - h_\epsilon \rangle_n| + |\langle w, h_\epsilon \rangle_n| \\ &\leq \frac{1}{n} \|w\|_1 \epsilon + |\langle w, h_\epsilon \rangle_n|\end{aligned}$$

Therefore

$$\sup_{g \in \mathcal{G}} |\langle w, h \rangle_n| \leq \frac{1}{n} \|w\|_1 \epsilon + \sup_{h_\epsilon \in \mathcal{G}_\epsilon} |\langle w, h_\epsilon \rangle_n|$$

The final goal is to show that $\sup_{h \in \mathcal{G}} |\langle w, h \rangle_n| \leq \text{constants} \cdot \sqrt{\frac{1}{n^{4/5}}}$.

Chaining.

Our goal is to give a precise bound for $\sup_{h \in \mathcal{G}} |\langle w, h \rangle|$.

We first restate the chaining theorem. Suppose $\|g\|_n \leq R$ for all $g \in \mathcal{G}$.

Suppose $\epsilon > \frac{1}{\sqrt{n}} \sigma \mathbf{c} \int_0^R \sqrt{\log N_2(t, \mathcal{G})} dt \vee R$. Then we have that

$$P\left(\sup_{g \in \mathcal{G}} \langle w, g \rangle_n \geq \epsilon\right) \leq 4 \exp\left(-\frac{n\epsilon^2}{\mathbf{c}R^2\sigma^2}\right)$$

In our case, $R \leq sLb$.

Restated, we have that, with probability at least $1 - \delta$,

$$\sup_{g \in \mathcal{G}} |\langle w, g \rangle_n| \leq \mathbf{c}sLb\sigma \sqrt{\frac{1}{n} \log \frac{4}{\delta}} + \left(\int_0^R \sqrt{\log N_2(t, \mathcal{G})} dt \vee R\right) \mathbf{c}\sigma \sqrt{\frac{1}{n}}$$

Now we evaluate the integral. Since $N_2(t, \mathcal{G}) \leq N_\infty(t, \mathcal{G})$, we know that $\sqrt{\log N_2(t, \mathcal{G})} \leq \sqrt{Cs^{1.5}bLt^{-1/4}}$.

$$\begin{aligned} \int_0^R \sqrt{\log N_2(t, \mathcal{G})} dt &\leq \sqrt{Cs^{1.5}bL} \int_0^R t^{-1/4} dt \\ &= \sqrt{Cs^{1.5}bL} \frac{4}{3} R^{3/4} \\ &= \sqrt{Cs^{1.5}bL} \mathbf{c}(sLb)^{3/4} \\ &\leq \mathbf{c}(sbL)^2 \end{aligned}$$

Coming back, we have, with probability at least $1 - \delta$,

$$\begin{aligned} \sup_{g \in \mathcal{G}} |\langle w, g \rangle| &\leq \mathbf{c}sLb\sigma \sqrt{\frac{1}{n} \log \frac{4}{\delta}} + \mathbf{c}(sLb)^2 \sigma \sqrt{\frac{1}{n}} \\ &\leq \mathbf{c}(sLb)^2 \sigma \sqrt{\frac{1}{n} \log \frac{4}{\delta}} \end{aligned}$$

Uniform convergence:

The goal is to bound $\sup_{f \in \mathcal{C}_L^s} \left| \|f_0 - f\|_n^2 - \|f_0 - f\|_P^2 \right|$.

Again, let \mathcal{G} denote the off-centered set of convex functions, that is, $\mathcal{G} \equiv \mathcal{C}_L^s - f_0$.

First, note that if $h \in \mathcal{G}$, then $\|h\|_\infty = \|f_0 - f\|_\infty \leq 4sbL$.

Because $\|h\|_n = \|h - h_\epsilon + h_\epsilon\|_n$, we have that

$$\begin{aligned} \|h_\epsilon\|_n - \|h - h_\epsilon\|_n &\leq \|h\|_n \leq \|h_\epsilon\|_n + \|h - h_\epsilon\|_n \\ \|h_\epsilon\|_n - \epsilon &\leq \|h\|_n \leq \|h_\epsilon\|_n + \epsilon \\ \|h_\epsilon\|_n^2 - 8\epsilon(sLb) &\leq \|h\|_n^2 \leq \|h_\epsilon\|_n^2 + 8\epsilon(sLb) \end{aligned}$$

where we used the fact that $\|h - h_\epsilon\|_n \leq \|h - h_\epsilon\|_\infty \leq \epsilon$ and $\|h_\epsilon\|_n \leq \|h_\epsilon\|_\infty \leq 4sbL$.

And likewise:

$$\|h_\epsilon\|_P^2 - 8\epsilon(sLb) \leq \|h\|_P^2 \leq \|h_\epsilon\|_P^2 + 8\epsilon(sLb)$$

Therefore,

$$\sup_{h \in \mathcal{G}} \left| \|h\|_n^2 - \|h\|_P^2 \right| \leq \sup_{h_\epsilon} \left| \|h_\epsilon\|_n^2 - \|h_\epsilon\|_P^2 \right| + \epsilon(16sLb)$$

Since $\|h_\epsilon\|_n^2 = \frac{1}{n} \sum_{i=1}^n h_\epsilon(X_i)^2$ is an average of bounded random variables, we have by Union Bound and Hoeffding Inequality that, with probability at most $1 - \delta$,

$$\begin{aligned} \sup_{h_\epsilon} \left| \|h_\epsilon\|_n^2 - \|h_\epsilon\|_P^2 \right| &\leq (8sLb)^2 \sqrt{\frac{1}{cn} \left(\log \frac{2}{\delta} + \log N_\infty(\epsilon, \mathcal{C}_L^s) \right)} \\ &\leq (8sLb)^2 \sqrt{\frac{1}{cn} \left(\log \frac{2}{\delta} + Cs^{1.5}Lb\epsilon^{-1/2} \right)} \end{aligned}$$

We will set $\epsilon = \frac{1}{n^{2/5}}(Cs^{0.5}Lb)^2$ and $\delta = \frac{1}{n}$. Therefore:

$$\begin{aligned} \sup_{h_\epsilon} \left| \|h_\epsilon\|_n^2 - \|h_\epsilon\|_P^2 \right| &\leq (8sLb)^2 \sqrt{\frac{1}{cn} \left(\log \frac{2}{\delta} + sn^{1/5} \right)} \\ &\leq (8Lb)^2 \sqrt{\frac{s^5}{cn^{4/5}} \log 2n} \end{aligned}$$

And

$$\begin{aligned} \sup_{h \in \mathcal{G}} \left| \|h\|_n^2 - \|h\|_P^2 \right| &\leq \sup_{h_\epsilon} \left| \|h_\epsilon\|_n^2 - \|h_\epsilon\|_P^2 \right| + \frac{1}{n^{2/5}} C^2 s^2 (Lb)^3 \\ &\leq (8Lb)^2 \sqrt{\frac{s^5}{cn^{4/5}} \log 2n} + (CLb)^2 \sqrt{\frac{s^4}{n^{4/5}}} \\ &\leq \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 2n} \end{aligned}$$

Concavity start.

We start from the very beginning, where \widehat{g}_k is defined as $\min_{g_k \in \mathcal{C}_L^1} \|y - \widehat{f} - g_k\|_n^2 + \lambda \|g_k\|_\infty$.

$$\begin{aligned} \|y - \widehat{f} - \widehat{g}_k\|_n^2 + \lambda \|\widehat{g}_k\|_\infty &\leq \|y - \widehat{f} - g_k^*\|_n^2 + \lambda \|g_k^*\|_\infty \\ \|y - \widehat{f} - \widehat{g}_k\|_n^2 &\leq \|y - \widehat{f} - g_k^*\|_n^2 + \lambda 2Lb \end{aligned}$$

$$\begin{aligned} \|f_0 - \widehat{f} - \widehat{g}_k + w\|_n^2 &\leq \|f_0 - \widehat{f} - g_k^* + w\|_n^2 + \lambda 2Lb \\ \|f_0 - \widehat{f} - \widehat{g}_k\|_n^2 - \|f_0 - \widehat{f} - g_k^*\|_n^2 &\leq 2\langle w, \widehat{g}_k - g_k^* \rangle_n + \lambda 2Lb \end{aligned}$$

Using the chaining result, setting $s = 1$, we have, with probability at least $1 - \delta$,

$$\|f_0 - \widehat{f} - \widehat{g}_k\|_n^2 - \|f_0 - \widehat{f} - g_k^*\|_n^2 \leq \mathbf{c}(Lb)^2 \sigma \sqrt{\frac{1}{n} \log \frac{4}{\delta}} + \lambda 2Lb$$

Setting $\delta = 1/n$, using *uniform convergence*, we have

$$\begin{aligned}\|f_0 - \hat{f} - \hat{g}_k\|_P^2 - \|f_0 - \hat{f} - g_k^*\|_P^2 &\leq \mathbf{c}(Lb)^2 \sigma \sqrt{\frac{1}{n} \log 4n + \lambda 2Lb} + \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 2n} \\ &\leq \mathbf{c}(Lb)^3 \sigma \sqrt{\frac{s^5}{n^{4/5}} \log 4n + \lambda 2Lb}\end{aligned}$$

Now comes **concavity removal version 2**.

The goal is to bound radius of approximation between $\|f_0 - \hat{f} - g_k^*\|_P^2$ and $\|f_0 - f^* - g_k^*\|_P^2$ and likewise for \hat{g}_k .

$$\begin{aligned}\|f_0 - \hat{f} - g_k^*\|_P^2 - \|f_0 - f^* - g_k^*\|_P^2 &\leq \|f_0 - \hat{f}\|_P^2 - \|f_0 - f^*\|_P^2 - 2\langle f_0 - \hat{f}, g_k^* \rangle + 2\langle f_0 - f^*, g_k^* \rangle \\ &\leq \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 4n + \lambda 2sLb} + 2|\langle \hat{f} - f^*, g_k^* \rangle| \\ &\leq \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 4n + \lambda 2sLb} + 2\|\hat{f} - f^*\|_P \|g_k^*\|_P \\ &\leq \mathbf{c}(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 4n + \lambda 2sLb} + \mathbf{c}2Lb \sqrt{(Lb)^3 \sqrt{\frac{s^5}{n^{4/5}} \log 4n + \lambda 2sLb}}\end{aligned}$$

The same bound likewise holds for \hat{g}_k .

2 2/15/2014

The two step procedure for variable selection:

1. Jointly convex approximation:

$$\{f_k^c\}_{k=1,\dots,p} = \arg \min_{f_k \in \mathcal{C}_x} \mathbb{E} \left(f(X) - \sum_{k=1}^p f_k(X_k) \right)^2$$

2. Separate concave post-processing, for each $k = 1, \dots, p$:

$$f_k^v = \arg \min_{f_k \in \mathcal{C}_v} \mathbb{E} \left(f(X) - \sum_{k' \neq k} f_{k'}^c(X_{k'}) - f_k(X_k) \right)^2$$

We consider as irrelevant any k where $f_k^c = 0$ and $f_k^v = 0$.

Claim: Let $f : C \rightarrow \mathbb{R}$ where $C = [0, 1]^p$. Suppose $p(x)$ is a positive density on C and that it satisfies the boundary-points condition.

Suppose $f(x)$ and $\nabla f(x)$ are continuous and for all k , that $\partial_{x_k} p(x | x_k)$ and $\partial_{x_k}^2 p(x | x_k)$ are also continuous as functions on C .

Then, the two step procedure is faithful, that is, $f_k^c = 0$ and $f_k^v = 0$ implies that $\partial_{x_k} f(x) = 0$.

Proof. We focus on a specific k and suppose that f_k^c and f_k^v are both 0.

We now invoke proposition 2.1 by letting $g(x) = f(x) - \sum_{k' \neq k} f_{k'}^c(x_{k'})$. It is easy to verify that derivative continuity requirements on g and $p(x | x_k)$ are all met. We can therefore conclude that $g_k^*(x) = \mathbb{E}[f(X) - \sum_{k' \neq k} f_{k'}^c(X_{k'}) | x_k] = 0$.

We now use corollary 2.1 by letting $\phi(x_{-k}) = \sum_{k' \neq k} f_{k'}^c(x_{k'})$. We therefore derive that $\partial_{x_k} f(x) = 0$ and thus prove faithfulness. □

2.1 Analysis

To prove the faithfulness of this procedure, we use the following extension of additive faithfulness.

Corollary 2.1. *Suppose $p(x)$ is a positive density on $[0, 1]^p$ and it satisfies the boundary-points condition.*

For any function $\phi(X_{-k})$ that does not depend on X_k :

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

We have that $f_k^* = 0 \Rightarrow \partial_{x_k} f(x) = 0$.

Proof. Identical to the proof of additive faithfulness. □

We also use the following form of the shape-constrained projection theorem.

Proposition 2.1. *Let $C \subset \mathbb{R}^p$ be a compact set and let $g : C \rightarrow \mathbb{R}$. Let $p(x)$ be a positive density on C and suppose $\mathbb{E}g(X) = 0$.*

Suppose that $\partial_{x_k} g(x)$, $\partial_{x_k} p(x | x_k)$, and $\partial_{x_k}^2 p(x | x_k)$ are all continuous as functions on C . Suppose that $\partial_{x_k}^2 g(x) \geq 0$.

Let $f_k^c(x_k) = \arg \min_{f_k \in \mathcal{C}_x} \mathbb{E} \left(g(X) - f_k(X_k) \right)^2$ and $f_k^v(x_k) = \arg \min_{f_k \in \mathcal{C}_v} \mathbb{E} \left(g(X) - f_k(X_k) \right)^2$ be the best convex and concave univariate approximation respectively.

Then, $f_k^c = 0$ and $f_k^v = 0$ iff $g_k^*(x_k) = \mathbb{E}[g(X) | x_k] = 0$.

Proof. First, we will establish that $g_k^*(x_k)$ is twice differentiable and that $\partial_{x_k}^2 g_k^*(x_k)$ is lower bounded.

$$\begin{aligned} g_k^*(x_k) &= \mathbb{E}[g(X) | x_k] \\ &= \int_{x-k} g(x) p(x | x_k) \\ \partial_{x_k}^2 g_k^*(x_k) &= \int_{x-k} g''(x) p(x | x_k) + 2g'(x) p'(x | x_k) + g(x) p''(x | x_k) dx_{-k} \end{aligned}$$

The first term $g''(x) p(x | x_k)$ is strictly positive. By assumption, the remaining terms are continuous and hence bounded on a compact set. $\partial_{x_k}^2 g_k^*(x_k)$ is therefore lower bounded.

Before proceeding, we also note that because $\mathbb{E}g(X) = 0$, it must be that $\mathbb{E}g_k^*(X_k) = 0$.

Now suppose $f_k^c = 0$ and $f_k^v = 0$. Let σ_k^2 denote $\mathbb{E}X_k^2$. Then

$$\arg \min_{c \in \mathbb{R}} \mathbb{E} \left(g(X) - c(X_k^2 - \sigma_k^2) \right)^2 = 0$$

Since optimal $c^* = \frac{\mathbb{E}[g(X)(X_k^2 - \sigma_k^2)]}{\mathbb{E}[X_k^2]}$, we know $\mathbb{E}[g(X)X_k^2] = \mathbb{E}[\mathbb{E}[g(X) | X_k]X_k^2] = 0$.

Because $\partial_{x_k}^2 g_k^*(x_k)$ is lower bounded, for large enough α , $g_k^*(x_k) + \alpha(x_k^2 - \sigma_k^2)$ has a non-negative second derivative and thus is convex. Then

$$\arg \min_{c \in \mathbb{R}} \mathbb{E} \left(g(X) - c(g_k^*(X_k) + \alpha(X_k^2 - \sigma_k^2)) \right)^2 = 0$$

Again, $c^* = \frac{\mathbb{E}[g(X)(g_k^*(X_k) + \alpha(X_k^2 - \sigma_k^2))]}{\mathbb{E}(g_k^*(X_k) + \alpha(X_k^2 - \sigma_k^2))^2} = 0$, so

$$\begin{aligned} \mathbb{E}[g(X)(g_k^*(X_k) + \alpha X_k^2)] &= \mathbb{E}[g(X)g_k^*(X_k)] \\ &= \mathbb{E}[\mathbb{E}[g(X) | X_k]g_k^*(X_k)] \\ &= \mathbb{E}g_k^*(X_k)^2 = 0 \end{aligned}$$

Therefore, $g_k^*(x_k) = 0$. □

3 2/7/2014

3.1 Best Shape-Constrained Projection

Even if a function f_0 is non-zero, it may be possible that the best convex approximation is zero. But, it cannot be that the best convex approximation and the best concave approximation are both zero.

Theorem 3.1. *Let $f_0 : C \rightarrow \mathbb{R}$ be a function with a bounded Hessian; let p be some distribution on C . Suppose*

$$\operatorname{argmin}_{f \in \mathcal{C}_x \cup \mathcal{C}_v} \mathbb{E}(f_0(x) - f(x))^2 = 0$$

then $f_0 = 0$ necessarily.

Proof. Suppose that the $\operatorname{argmin}_{f \in \mathcal{C}_x \cup \mathcal{C}_v} \mathbb{E}(f_0(x) - f(x))^2 = 0$, then it must be that

$$\operatorname{argmin}_{c \in \mathbb{R}} \mathbb{E} \left(f_0(x) - cx^\top x \right)^2 = 0$$

Since the optimal c in the above optimization is $c^* = \frac{\mathbb{E}[f_0(x)x^\top x]}{\mathbb{E}(x^\top x)^2}$, we have that $\mathbb{E}[f_0(x)x^\top x] = 0$.

We know that there exists a convex function f' such that $f' = f_0 + rx^\top c$ for some large $r > 0$.

By assumption and convexity of f' , it must be that

$$\operatorname{argmin}_{c \in \mathbb{R}} \mathbb{E} \left(f_0(x) - cf'(x) \right)^2 = 0$$

We know by similar argument then that $\mathbb{E}[f_0(x)f'(x)] = 0$.

However, $\mathbb{E}[f_0(x)f'(x)] = \mathbb{E}[f_0(x)^2 + rf_0(x)x^\top x] = \mathbb{E}[f_0(x)^2]$. This is a contradiction. □

Thus, if we tried both convex and concave projection and the solution is zero, then we can be sure that the original function is identically zero.

The same argument applies if we consider the best single component approximation:

$$\operatorname{argmin}_{f_k \in \mathcal{C}_x \cup \mathcal{C}_v} \mathbb{E} \left(f_0(x) - f_k(x_k) \right)^2$$

Because $\mathbb{E}[f_0(x)x_k^2] = \mathbb{E} \left[\mathbb{E}[f_0(x) | x_k] x_k^2 \right] = \mathbb{E}[f_k^*(x_k)x_k^2]$. Thus, we can apply the same argument and use $f_k^*(x_k)$ where needed.

4 1/23/2014

4.1 Additive Faithfulness Case Study with Quadratic Function and Gaussian Distribution

We consider a quadratic function $f(x) = x^\top Hx + c^\top x$ and a Gaussian distribution $X \sim N(0, \Sigma)$.

We then have a closed form for the additive approximation.

- If f does not depend on x_j , then $f_j^*(x_j) = 0$.
- If f depends on x_j , then, letting H_j be the j -th row of H and Σ_j be the j -th row of Σ :

$$f_j^*(x_j) = H_j^\top \Sigma_j \frac{1}{\Sigma_{jj}} x_j^2 + c_j x_j$$

Let us assume $\Sigma_{jj} = 1$ for all j and that $c = 0$ for convenience. We then have two direct corollaries:

Corollary: We have additive convexity if and only if $\text{diag}(H\Sigma) \geq 0$.

Corollary: We have additive faithfulness if and only if $\text{diag}(H\Sigma) \neq 0$.

As an example where additive convexity and additive faithfulness are violated. Let $H = [1, 2; 2, 5]$ and $\Sigma = [1, -c; -c, 1]$. For $c = 0.5$, additive faithfulness is violated; for $c > 0.5$, additive convexity is violated.

Proof. We will show that the $f_j^*(x_j)$'s, so described, satisfy the KKT stationarity equations.

$$f_j^*(x_j) = \mathbb{E}[f(x) - \sum_{k \neq j} f_k^*(x_k) | x_j] \quad \text{for all } j$$

To prove this, we use the following conditional mean and conditional covariance property of the multivariate Gaussian distribution.

$$\begin{aligned} \mathbb{E}[x_k | x_j] &= \Sigma_{jk} \Sigma_{jj}^{-1} x_j \\ \mathbb{E}[x_k x_{k'} | x_j] &\text{ is some constant for all } x_j \\ \mathbb{E}[x_k^2 | x_j] &\text{ is some constant for all } x_j \end{aligned}$$

□

Why can Gaussian distribution violate additive faithfulness? Because $\frac{\partial p(x_{-j} | x_j)}{\partial x_j}$ is always large for some values of x_{-j} .

5 1/2/2014

5.1 Convex-minus-Quadratic Estimation

Instead of estimating a convex-plus-concave function, it is theoretically sufficient to estimate a convex-minus-quadratic function.

Theorem 5.1. *Any function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ with a bounded Hessian can be decomposed as $h(x) = f(x) - q(x)$ where $f(x)$ is convex and $q(x)$ is $cx^\top x$ for some $c \geq 0$.*

Proof. For large enough c , the Hessian of $h + c$ is a positive semidefinite. The function f is thus convex. \square

Convex-minus-quadratic functions are faster to learn. The optimization program has about twice as few variables. It is also possibly easier to analyze.

5.1.1 Implementation

For additive modeling, I use the following optimization program. There are two parameters λ and L .

$$\begin{aligned} \min_{f_1, \dots, f_p} \quad & \frac{1}{n} \sum_i \left(y_i - \sum_j f_j(x_{ij}) + c_j x_{ij}^2 \right)^2 + \lambda \sum_j \|\partial f_j - 2c_j x_j^2\|_\infty \\ \text{s.t.} \quad & f_j \text{ is convex} \\ & c_j \geq 0 \quad \text{and} \quad c_j \leq L \end{aligned}$$

The second L parameter is necessary. A similar parameter is required in convex-plus-concave estimation as well, which I will explain in the next section. L can be interpreted as a lower bound on the second derivative of the estimated regression function. We set $L = 200$ in experiments.

Demonstration: We estimate a one-dimensional function so that we can visualize the behavior of convex-minus-quadratic functions. We set $\lambda = 0$. The result is in figure 5.1.1.

Here, $n = 300$, the SNR is about 0.6 (high noise).

5.2 Tuning Parameters for Convex-plus-Concave Estimation

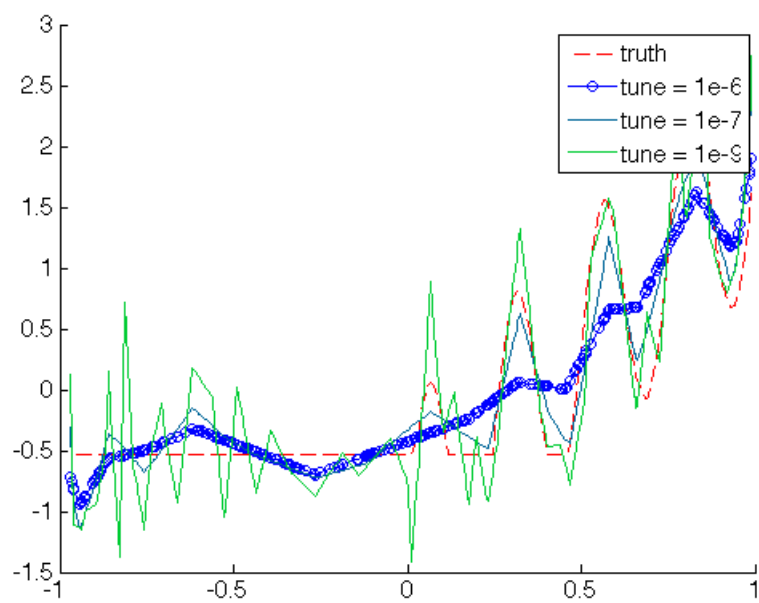
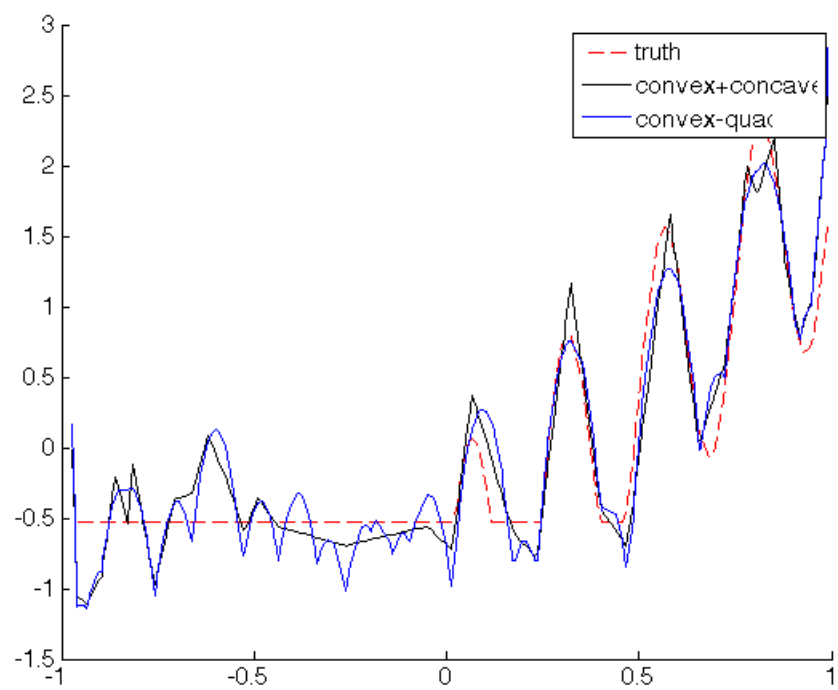
Convex-plus-concave estimation is not free of tuning parameters. Arbitrary sum of convex and concave functions can represent any function with a bounded Hessian.

In Minhua's original SCCAM implementation, the objective is augmented with a square-penalty on the value of the gradients of f and g , the convex and the concave functions. The penalty looks like $10^{-6} \|\text{gradient}\|_2^2$. This penalty is not just for numerical stability; its magnitude affects the estimation.

Demonstration:

The experimental set-up is the same as before. The result is in figure 5.2

As one can see, when the tuning parameter (the coefficient for the $\|\text{gradient}\|_2^2$ penalty) is too small, the estimated function is fitting noise. When the tuning parameter is too large, the estimated function is too smooth.



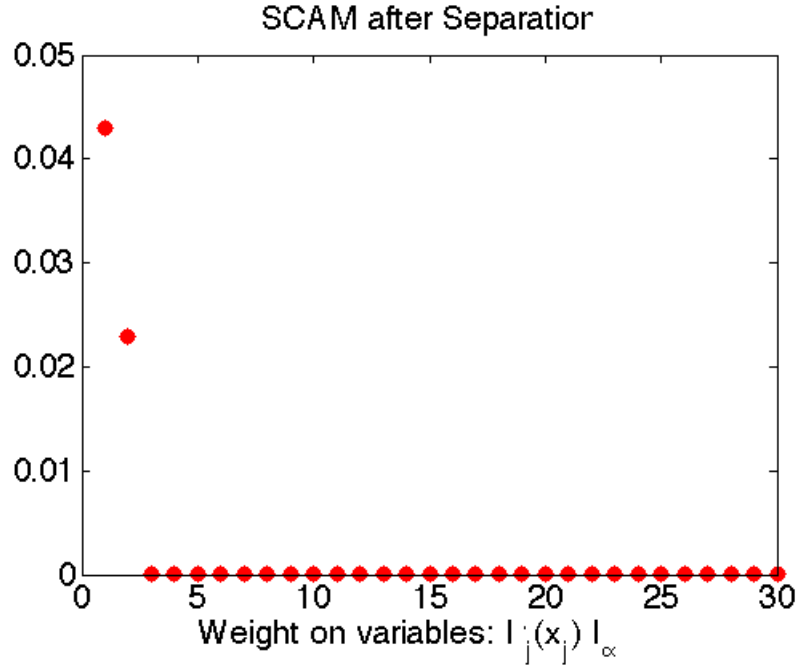
5.3 Convex-Concave Separation

We study the problem of variable selection on a convex-plus-concave function h . The approach has two steps:

1. Learn **sparse non-additive** convex function f and concave function g such that $h = f + g$. This is the separation stage.
2. Apply sparse additive model on f and g .

Preliminary results are favorable. In our experiment, we let $p = 30, n = 300$. The true function is $h(x) = 2x_1x_2$ and uses only the first two variables. The SNR is 4. h is an example of a function that cannot be consistently estimated by additive modeling.

We can indeed identify the correct variables by applying additive model on f after separation.



Without separation however, the additive model indeed fails.

