

Model-assisted survey regression estimation with the lasso ¹

F. Jay Breidt
Colorado State University

Opening Workshop on Computational Methods in Social Sciences
SAMSI August 2013

This research was supported in part by the US National Science Foundation (SES0922142).

Model-assisted survey regression estimation with the lasso

2

- (Model-assisted survey regression estimation):
 - joint work with Jean Opsomer, Colorado State, and various other colleagues
 - bringing in nonparametric and semiparametric regression methods into classical survey statistics
- (with the lasso):
 - joint work with former PhD student Kelly McConville, Whitman College; Thomas Lee, UC-Davis; and Gretchen Moisen, US Forest Service

- Finite population $U = \{1, 2, \dots, N\}$
- Response variable $y_k, k \in U$
 - these are just unknown real numbers, with no probability structure
- Goal: estimate finite population total

$$t_y = \sum_{k \in U} y_k$$

- Draw probability sample $s \subset U$ with $\Pr[k \in s] = \pi_k > 0$

- Sample membership indicator $I_k = 1$ if $k \in s$, $I_k = 0$ otherwise

$E[I_k] = \pi_k$, averaging over all possible samples

- Use this repeated-sampling probability structure for statistical inference
- Unbiased Horvitz-Thompson estimator of t_y is

$$\text{HT}(y_k) = \sum_{k \in s} \frac{y_k}{\pi_k} = \sum_{k \in U} \frac{y_k I_k}{\pi_k}$$

- Depends on covariance structure of $\{I_k\}_{k \in U}$:

$$\begin{aligned} \text{Var} \left(\sum_{k \in U} y_k \frac{I_k}{\pi_k} \right) &= \sum_{k, l \in U} \text{Cov} (I_k, I_l) \frac{y_k y_l}{\pi_k \pi_l} \\ &= \sum_{k, l \in U} \Delta_{kl} \frac{y_k y_l}{\pi_k \pi_l} \end{aligned}$$

where $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$ and $\pi_{kl} = \Pr [I_k = 1, I_l = 1]$

Variance estimator for the Horvitz-Thompson estimator 6

- Provided $\pi_{kl} > 0$ for all $k, l \in U$,

$$\widehat{V}(\text{HT}) = \sum_{k,l \in U} \Delta_{kl} \frac{y_k y_l I_k I_l}{\pi_k \pi_l \pi_{kl}}$$

is unbiased for

$$\text{Var}(\text{HT}) = \sum_{k,l \in U} \Delta_{kl} \frac{y_k y_l}{\pi_k \pi_l}$$

- Under very mild **design** conditions,
 - $\widehat{V}(\text{HT})$ is consistent for $\text{Var}(\text{HT})$
 - HT is asymptotically normal and confidence intervals can be based on this fact in moderate to large samples

$$\left\{ \widehat{V}(\text{HT}) \right\}^{-1/2} \left\{ \text{HT}(y_k) - t_y \right\} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$$

- Involves only “finite second moments” of $\{y_k\}$,
 - not distributional assumptions
 - not dependence assumptions
- Good design makes normal approximations better

- Suppose we have auxiliary information vector \mathbf{x}_k , $k \in U$
- Also have a “method” $m(\cdot)$ for predicting y_k from \mathbf{x}_k :

$$y_k \simeq m(\mathbf{x}_k)$$

- method $m(\cdot)$ does not depend on the sample
- e.g., inflation-adjust an old census value
- Unbiased difference estimator of t_y is then

$$\text{Diff}_m(y_k) = \sum_{k \in U} m(\mathbf{x}_k) + \sum_{k \in s} \frac{y_k - m(\mathbf{x}_k)}{\pi_k}$$

$$\begin{aligned} & \text{Var} \left(\sum_{k \in U} m(\mathbf{x}_k) + \sum_{k \in U} (y_k - m(\mathbf{x}_k)) \frac{I_k}{\pi_k} \right) \\ &= \sum_{k, \ell \in U} \Delta_{k\ell} \frac{y_k - m(\mathbf{x}_k)}{\pi_k} \frac{y_\ell - m(\mathbf{x}_\ell)}{\pi_\ell} \end{aligned}$$

- Compare to Horvitz-Thompson estimator:

$$\text{Var} \left(\sum_{k \in U} y_k \frac{I_k}{\pi_k} \right) = \sum_{k, \ell \in U} \Delta_{k\ell} \frac{y_k}{\pi_k} \frac{y_\ell}{\pi_\ell}$$

- Provided $\pi_{kl} > 0$ for all $k, l \in U$,

$$\widehat{V}(\text{Diff}_m) = \sum_{k, l \in U} \Delta_{kl} \frac{y_k - m(\mathbf{x}_k)}{\pi_k} \frac{y_l - m(\mathbf{x}_l)}{\pi_l} \frac{I_k I_l}{\pi_{kl}}$$

is unbiased for

$$\text{Var}(\text{Diff}_m) = \sum_{k, l \in U} \Delta_{kl} \frac{y_k - m(\mathbf{x}_k)}{\pi_k} \frac{y_l - m(\mathbf{x}_l)}{\pi_l}$$

- Inherits asymptotic normality from HT under mild conditions.

- Difference estimator is **exactly** unbiased, regardless of the quality of the method $m(\cdot)$

$$\text{Diff}_m(y_k) = \sum_{k \in U} m(\mathbf{x}_k) + \text{HT}(y_k - m(\mathbf{x}_k))$$

- Has smaller variance than $\text{HT}(y_k)$ provided “residuals”

$$y_k - m(\mathbf{x}_k)$$

have smaller variation than “raw values” y_k

- Have an **exactly** unbiased variance estimator
- Results require that $m(\cdot)$ does not depend on the sample

- Difference estimator requires method $m(\cdot)$ independent of the sample
- Model-assisted estimator introduces a **working model**

$$y_k = \mu(\mathbf{x}_k) + \epsilon_k$$

- If the entire population were observed, use a standard statistical method to estimate $\mu(\cdot)$ by $m_N(\cdot)$ (independent of sample)

- Since only a sample is observed, estimate $m_N(\cdot)$ by $\hat{m}_N(\cdot)$
 - not independent of the sample
- Plug $\hat{m}_N(\cdot)$ into the difference estimator form:

$$\text{MA}(y_k) = \sum_{k \in U} \hat{m}_N(\mathbf{x}_k) + \sum_{k \in s} \frac{y_k - \hat{m}_N(\mathbf{x}_k)}{\pi_k}$$

- Plug $\hat{m}_N(\cdot)$ into the variance estimator:

$$\hat{V}(\text{MA}) = \sum_{k, l \in U} \Delta_{kl} \frac{y_k - \hat{m}_N(\mathbf{x}_k)}{\pi_k} \frac{y_l - \hat{m}_N(\mathbf{x}_l)}{\pi_l} \frac{I_k I_l}{\pi_{kl}}$$

- Working model is heteroskedastic multiple regression:

$$y_k = \mu(\mathbf{x}_k) + \epsilon_k = \mathbf{x}'_k \boldsymbol{\beta} + \epsilon_k, \quad \epsilon_k \sim (0, \sigma_k^2)$$

- If the entire population were observed, use weighted least squares:

$$m_N(\mathbf{x}_k) = \mathbf{x}'_k \mathbf{B}_N = \mathbf{x}'_k \left(\sum_{j \in U} \frac{\mathbf{x}_j \mathbf{x}'_j}{\sigma_j^2} \right)^{-1} \sum_{j \in U} \frac{\mathbf{x}_j y_j}{\sigma_j^2}$$

- Estimate finite population fit from the observed sample:

$$\hat{m}_N(\mathbf{x}_k) = \mathbf{x}'_k \widehat{\mathbf{B}}_N = \mathbf{x}'_k \left(\sum_{j \in s} \frac{\mathbf{x}_j \mathbf{x}'_j}{\pi_j \sigma_j^2} \right)^{-1} \sum_{j \in s} \frac{\mathbf{x}_j y_j}{\pi_j \sigma_j^2}$$

- $\widehat{\mathbf{B}}_N$ is asymptotically design unbiased and consistent for \mathbf{B}_N regardless of the quality of the working model specification

- Plug into model-assisted form:

$$\text{GREG}(y_k) = \sum_{k \in U} \mathbf{x}'_k \widehat{\mathbf{B}}_N + \sum_{k \in s} \frac{y_k - \mathbf{x}'_k \widehat{\mathbf{B}}_N}{\pi_k}$$

- Plug into the variance estimator:

$$\widehat{V}(\text{GREG}) = \sum_{k, \ell \in U} \Delta_{k\ell} \frac{y_k - \mathbf{x}'_k \widehat{\mathbf{B}}_N}{\pi_k} \frac{y_\ell - \mathbf{x}'_\ell \widehat{\mathbf{B}}_N}{\pi_\ell} \frac{I_k I_\ell}{\pi_{k\ell}}$$

- Classical survey ratio estimator: x_k is scalar, model is heteroskedastic regression through the origin
- Classical survey regression estimator: x_k is scalar, model is homoskedastic simple linear regression
- Post-stratification estimator: x_k is vector of indicators for categorical covariate
- . . .

- Write GREG as

$$\begin{aligned}
 \text{GREG}(y_k) &= \sum_{k \in U} \mathbf{x}'_k \widehat{\mathbf{B}}_N + \sum_{k \in s} \frac{y_k - \mathbf{x}'_k \widehat{\mathbf{B}}_N}{\pi_k} \\
 &= \sum_{k \in U} \mathbf{x}'_k \mathbf{B}_N + \sum_{k \in U} \frac{(y_k - \mathbf{x}'_k \mathbf{B}_N) I_k}{\pi_k} \\
 &\quad + \left(\widehat{\mathbf{B}}_N - \mathbf{B}_N \right)' \sum_{k \in U} \mathbf{x}_k \left(1 - \frac{I_k}{\pi_k} \right) \\
 &= \text{Diff}_{m_N}(y_k) + (\text{smaller-order term})
 \end{aligned}$$

- Asymptotically unbiased (and mean square consistent), regardless of the quality of the working model $\mu(\cdot)$
- Variance is asymptotically equivalent to

$$\begin{aligned} \text{Var} & \left(\sum_{k \in U} \mathbf{x}'_k \mathbf{B}_N + \sum_{k \in U} (y_k - \mathbf{x}'_k \mathbf{B}_N) \frac{I_k}{\pi_k} \right) \\ & = \sum_{k, \ell \in U} \Delta_{k\ell} \frac{y_k - \mathbf{x}'_k \mathbf{B}_N y_\ell - \mathbf{x}'_k \mathbf{B}_N}{\pi_k \pi_\ell} \end{aligned}$$

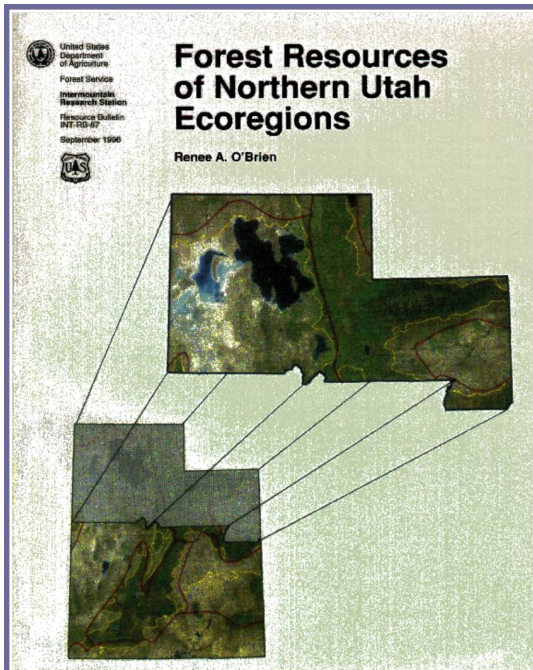
- Smaller asymptotic variance than HT(y_k) provided residuals $y_k - \mathbf{x}'_k \mathbf{B}_N$ have less variation than raw values y_k

- GREG can also be written in weighted form:

$$\begin{aligned}
 \text{GREG}(y_k) &= \sum_{k \in s} \frac{y_k - \mathbf{x}'_k \widehat{\mathbf{B}}_N}{\pi_k} + \sum_{k \in U} \mathbf{x}'_k \widehat{\mathbf{B}}_N \\
 &= \sum_{k \in s} \left\{ \frac{1}{\pi_k} + (\mathbf{t}\mathbf{x} - \text{HT}(\mathbf{x}_k))' \left(\sum_{k \in s} \frac{\mathbf{x}_k \mathbf{x}'_k}{\pi_k} \right)^{-1} \frac{\mathbf{x}_k}{\pi_k} \right\} y_k \\
 &= \sum_{k \in s} \omega_{ks} y_k
 \end{aligned}$$

- The GREG weights $\{\omega_{ks}\}$ do not depend on y and can be applied to any response variable

- Estimates required for forest area, wood volume, growth, mortality, . . .
- By region, species and other classifications



- Note that the weights $\{\omega_{ks}\}$ are calibrated to the \mathbf{X} -totals:

$$\begin{aligned}\text{GREG}(\mathbf{x}'_k) &= \sum_{k \in s} \left\{ \frac{1}{\pi_k} + (t\mathbf{x} - \text{HT}(\mathbf{x}_k))' \left(\sum_{k \in s} \frac{\mathbf{x}_k \mathbf{x}'_k}{\pi_k} \right)^{-1} \frac{\mathbf{x}_k}{\pi_k} \right\} \mathbf{x}'_k \\ &= \text{HT}(\mathbf{x}'_k) + (t\mathbf{x} - \text{HT}(\mathbf{x}_k))' \left(\sum_{k \in s} \frac{\mathbf{x}_k \mathbf{x}'_k}{\pi_k} \right)^{-1} \sum_{k \in s} \frac{\mathbf{x}_k \mathbf{x}'_k}{\pi_k} \\ &= t'\mathbf{x}\end{aligned}$$

- GREG will be very efficient if y_k is approximately a linear combination of \mathbf{x}_k

- $\text{GREG}(\mathbf{x}'_k) = t'_k \mathbf{x}$
- Calibration reproduces known population information
 - internal consistency across statistical system
 - reassuring for users, logistically convenient for agencies
- GREG weight adjustments may be large
 - extreme weights, negative weights are possible
 - many methods developed to trim or stabilize weights, including **ridge calibration** (Rao and Singh 1997, *ASA Proc.*; Park and Fuller 2005, *Survey Methodology*; Montanari and Ranalli 2009, Stats Canada workshop.)

- Keep the calibration, change the metric: Deville and Särndal (1992, *JASA*) minimize $d\left(\pi_k^{-1}, \omega_{ks}\right)$ subject to
$$\sum_{k \in s} \omega_{ks} \mathbf{x}_k = t \mathbf{x}$$
- Specify the working model more flexibly:
 - Local polynomial regression (Breidt and Opsomer 2000)
 - Neural nets (Montanari and Ranalli 2005)
 - Penalized splines (Breidt, Claeskens, Opsomer 2005); Regression splines (Goga 2005)
 - Generalized additive models (Opsomer, Breidt, Moisen, and Kauer-
mann 2007); Nonparametric additive models (Wang and Wang
2011)

- Write down your favorite $m_N(\cdot)$ you would use if the entire population were observed
- Create survey-weighted version, $\hat{m}_N(\cdot)$
- Plug in and write model-assisted estimator as

$$\begin{aligned} \text{MA}(y_k) &= \sum_{k \in U} \hat{m}_N(\mathbf{x}_k) + \sum_{k \in s} \frac{y_k - \hat{m}_N(\mathbf{x}_k)}{\pi_k} \\ &= \text{Diff}_{m_N}(y_k) + (\text{smaller-order term}) \end{aligned}$$

- Breidt and Opsomer (2000), *Ann. Stat.*
- Working model: $\mu(\cdot)$ is a smooth function of scalar x
- Estimate $\mu(\cdot)$ via **local polynomial regression**:

$$m_N(x_i) = (1, 0, \dots, 0) (\mathbf{X}'_{Ui} \mathbf{W}_{Ui} \mathbf{X}_{Ui})^{-1} \mathbf{X}'_{Ui} \mathbf{W}_{Ui} \mathbf{y}_U$$

where

$$\mathbf{X}_{Ui} = \left[1 \quad x_j - x_i \quad \cdots \quad (x_j - x_i)^q \right]_{j \in U}$$

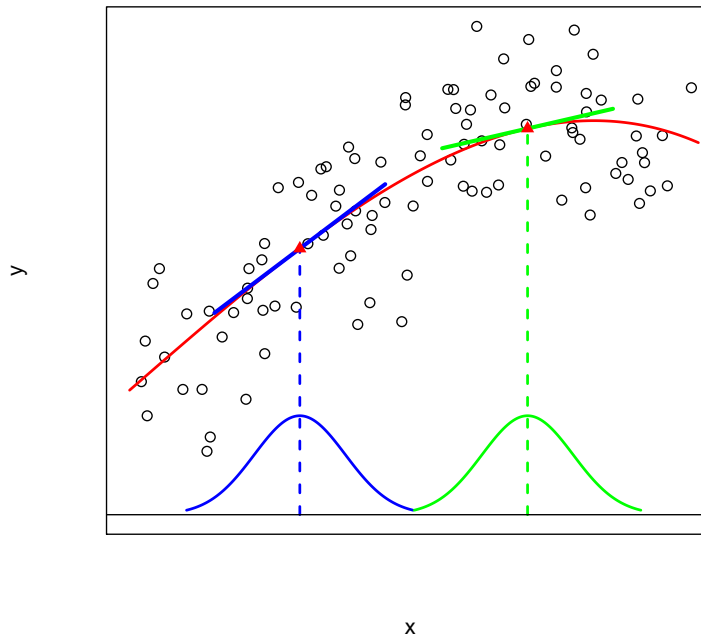
and

$$\mathbf{W}_{Ui} = \text{diag} \left\{ \frac{1}{h} K \left(\frac{x_j - x_i}{h} \right) \right\}_{j \in U}$$

Example: Local linear regression

27

- Fit line locally (as defined by kernel weights) and read off the local intercept



- Estimate $m_N(\cdot)$ using survey weights to get $\hat{m}_N(\cdot)$, plug in to model-assisted form

$$\begin{aligned}
 \text{LPR}(y_k) &= \sum_{k \in U} \hat{m}_N(\mathbf{x}_k) + \sum_{k \in s} \frac{y_k - \hat{m}_N(\mathbf{x}_k)}{\pi_k} \\
 &= \sum_{k \in U} m_N(\mathbf{x}_k) + \sum_{k \in U} \frac{(y_k - m_N(\mathbf{x}_k)) I_k}{\pi_k} \\
 &\quad + \sum_{k \in U} (\hat{m}_N(\mathbf{x}_k) - m_N(\mathbf{x}_k)) \left(1 - \frac{I_k}{\pi_k} \right) \\
 &= \text{Diff}_{m_N}(y_k) + (\text{smaller-order term??})
 \end{aligned}$$

- Consider sequence of finite populations with $n_N \rightarrow \infty$ as $N \rightarrow \infty$
- Smoothing assumptions:
 - kernel $K(\cdot)$ is symmetric, continuous, and compactly supported
 - bandwidth $h_N \rightarrow 0$ and $Nh_N^2 \rightarrow \infty$
- Design assumptions:
 - $\min_{i \in U_N} \pi_i \geq \lambda > 0$ and $\min_{i,j \in U_N} \pi_{ij} \geq \lambda^* > 0$
 - limited dependence:

$$\max_{i,j \in U_N: i \neq j} |\pi_{ij} - \pi_i \pi_j| = O(n_N^{-1})$$

$$\max_{(i,j,k,\ell) \text{ distinct}} |\mathbb{E}[(I_i - \pi_i)(I_j - \pi_j)(I_k - \pi_k)(I_\ell - \pi_\ell)]| = O(N^{-2})$$

$$\max_{(i,j,k,\ell) \text{ distinct}} |\mathbb{E}[(I_i I_j - \pi_{ij})(I_k I_\ell - \pi_{k\ell})]| = o(1)$$

$$\max_{(i,j,k) \text{ distinct}} |\mathbb{E}[(I_i - \pi_i)^2(I_j - \pi_j)(I_k - \pi_k)]| = O(n_N^{-1})$$

Under the above asymptotic framework,

- LPR estimator is mean square consistent for t_y
- Variance is asymptotically equivalent to

$$\sum_{k, \ell \in U} \Delta_{k\ell} \frac{y_k - m_N(\mathbf{x}_k)}{\pi_k} \frac{y_\ell - m_N(\mathbf{x}_\ell)}{\pi_\ell}$$

- Standard plug-in variance estimator is consistent
- Smaller asymptotic variance than HT provided residuals $y_k - m_N(\mathbf{x}_k)$ have less variation than raw values y_k

- $\text{LPR}(y_i) = \sum_{i \in s} \omega_{is} y_i$ with weights independent of y
- For q th order local polynomial, weights are calibrated to powers of x :

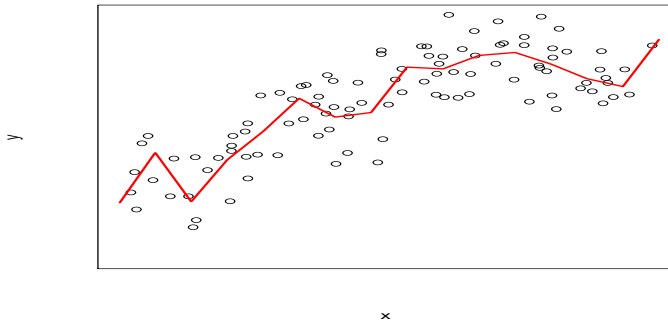
$$\sum_{i \in s} \omega_{is} x_i^\ell = \sum_{i \in U_N} x_i^\ell \quad (\ell = 0, 1, \dots, q)$$

- LPR will be particularly effective if y is approximately a q th order polynomial in x

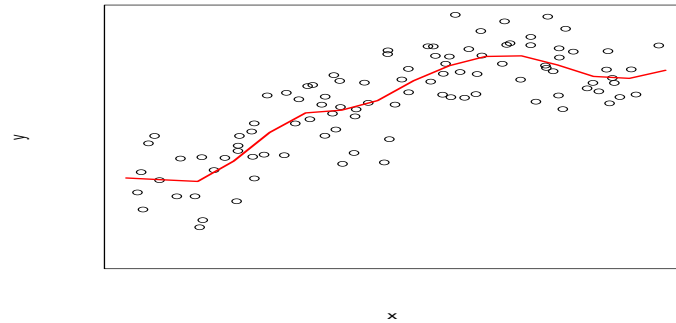
Example: Penalized spline regression

- K fixed, known knots and K basis functions
- Penalty parameter λ^2 determines degrees of freedom of the smooth

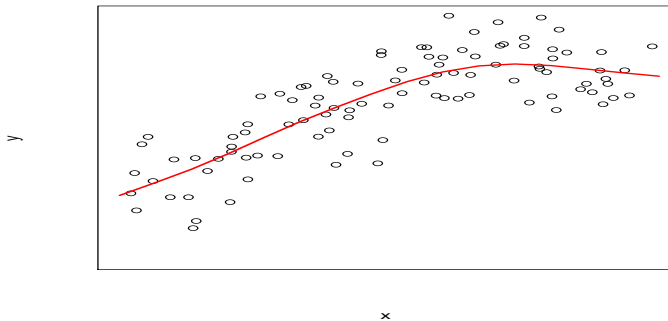
$\lambda^2 = 0$ and $df = 16$



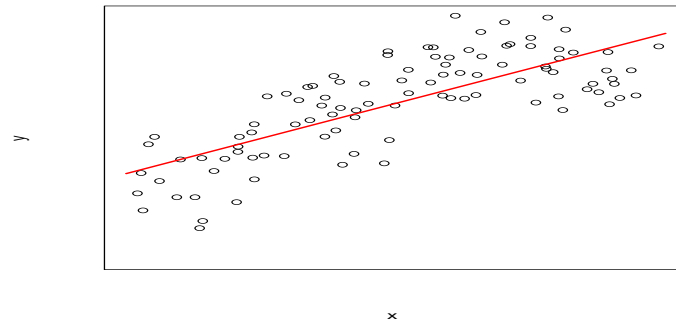
$\lambda^2 = 0.0083$ and $df = 8$



$\lambda^2 = 0.252$ and $df = 4$



$\lambda^2 = 1000$ and $df = 2$



- Breidt, Claeskens and Opsomer (2005, *Biometrika*)
- Choose penalty λ^2 to give specified degrees of freedom
- Formulate p-spline as linear mixed model (LMM)
- Write down LMM fit for the entire finite population:

$$[m_N(\mathbf{c}_k)]_{k \in U} = \mathbf{C}(\mathbf{C}'\mathbf{C} + \mathbf{\Lambda})^{-1}\mathbf{C}'\mathbf{y}_U$$

with $\mathbf{\Lambda} = \text{blockdiag}(\mathbf{0}, \lambda^2\mathbf{I}_K)$

- Estimate $m_N(\cdot)$ using survey weights to get $\hat{m}_N(\cdot)$, plug in to model-assisted form

- Under standard asymptotic framework with K fixed,
 - p-spline estimator is mean square consistent for t_y
 - variance is asymptotically equivalent to

$$\sum_{k, \ell \in U} \Delta_{k\ell} \frac{y_k - m_N(\mathbf{c}_k)}{\pi_k} \frac{y_\ell - m_N(\mathbf{c}_\ell)}{\pi_\ell}$$

- standard plug-in variance estimator is consistent
- smaller asymptotic variance than HT provided residuals $y_k - m_N(\mathbf{c}_k)$ have less variation than raw values y_k
- McConville and Breidt (2013, *J. Nonpar. Stat.*) prove above results with $K \rightarrow \infty$

- LMM model-assisted estimator (including p-spline) can be written

$$\begin{aligned}
 \text{LMM}(y_k) &= \sum_{k \in s} \frac{y_k - \mathbf{c}'_k \hat{\mathbf{B}}}{\pi_k} + \sum_{k \in U} \mathbf{c}'_k \hat{\mathbf{B}} \\
 &= \sum_{k \in s} \left\{ \frac{1}{\pi_k} + (t_{\mathbf{c}} - \text{HT}(\mathbf{c}_k))' \left(\sum_{k \in s} \frac{\mathbf{c}_k \mathbf{c}'_k}{\pi_k} + \mathbf{\Lambda} \right)^{-1} \frac{\mathbf{c}_k}{\pi_k} \right\} y_k \\
 &= \sum_{k \in s} \omega_{ks} y_k
 \end{aligned}$$

- Calibrated on \mathbf{X} , $\text{LMM}(\mathbf{x}'_k) = t'_{\mathbf{x}}$, but not on \mathbf{Z} , $\text{LMM}(\mathbf{z}'_k) \neq t'_{\mathbf{z}}$, due to the penalization

-
- Both LPR and p-splines have good robustness properties:
 - comparable efficiency to GREG when parametric working model is correct
 - better efficiency when parametric working model is incorrect
 - better-behaved weights than GREG (e.g., almost never negative)
 - P-splines extend much more readily than kernels to semi-parametric models:
 - additional \mathbf{X} variables, continuous or categorical, with calibration
 - additional \mathbf{Z} variables, continuous or categorical, without calibration

- Both LPR and p-splines had linear structure, allowing calibrated y -independent weights
- Nonlinear methods are typically uncalibrated, but...
- Wu and Sitter (2001) *JASA*: regress y_k on $\widehat{m}_N(\mathbf{x}_k)$:

$$\begin{aligned} \text{WS}(y_k) &= \sum_{k \in s} \frac{y_k - \widehat{m}_N(\mathbf{x}_k) \widehat{R}}{\pi_k} + \sum_{k \in U} \widehat{m}_N(\mathbf{x}_k) \widehat{R} \\ &= \sum_{k \in s} \left\{ \frac{1}{\pi_k} + \left(\sum_{k \in U} \widehat{m}_N(\mathbf{x}_k) - \sum_{k \in s} \frac{\widehat{m}_N(\mathbf{x}_k)}{\pi_k} \right) \frac{\widehat{m}_N(\mathbf{x}_k) / \pi_k}{\sum_{k \in s} \widehat{m}_N(\mathbf{x}_k)^2 / \pi_k} \right\} y_k \end{aligned}$$

yielding weights, but these depend on y

- Conducted by United States Forest Service
- Nationwide network of 400,000 sample plots, visited every five years in rotating panels
- **Goal:** national and regional estimates of t_y or \bar{y}_U for various y : forest area, wood volume, growth, mortality,...
 - **expensive** data y are field-collected or manually interpreted from aerial photography
 - interest in using **cheap** auxiliary data x from remote sensing and other spatial data sources

Auxiliary information for all $k \in U$ in real applications? 39

- Required auxiliary info comes from
 - “wall-to-wall” remote sensing, like satellite imagery
 - digital elevation models
 - geographic information system data products
- Working model is $\mu(\mathbf{x}_k) = \mathbf{x}'_k \boldsymbol{\beta}$, as with GREG:
 - but auxiliary variables may be highly correlated
 - may have poor predictive ability for response variables
- Consider using model selection methods

- Tibshirani (1996): Least absolute shrinkage and selection operator (lasso)

$$\mathbf{B}_N^{(L)} = \arg \min_{\boldsymbol{\beta}} (\mathbf{Y}_U - \mathbf{X}_U \boldsymbol{\beta})^T (\mathbf{Y}_U - \mathbf{X}_U \boldsymbol{\beta}) + \lambda_N \sum_{j=1}^p |\beta_j|$$

- simultaneously performs model selection and estimation by shrinking unnecessary coefficients to zero

- Construct survey-weighted version:

$$\widehat{\mathbf{B}}_N^{(L)} = \arg \min_{\boldsymbol{\beta}} (\mathbf{Y}_s - \mathbf{X}_s \boldsymbol{\beta})^T \boldsymbol{\Pi}_s^{-1} (\mathbf{Y}_s - \mathbf{X}_s \boldsymbol{\beta}) + \lambda_N \sum_{j=1}^p |\beta_j|$$

- Plug into model-assisted form:

$$\text{LASSO}(y_k) = \sum_{k \in U} \mathbf{x}'_k \widehat{\mathbf{B}}_N^{(L)} + \sum_{k \in s} \frac{y_k - \mathbf{x}'_k \widehat{\mathbf{B}}_N^{(L)}}{\pi_k}$$

- Assume conditions under which GREG is consistent and asymptotically equivalent to the difference estimator,

$$\sum_{k \in U} \mathbf{x}'_k \mathbf{B}_N + \sum_{k \in s} \frac{y_k - \mathbf{x}'_k \mathbf{B}_N}{\pi_k}, \text{ where } \mathbf{B}_N = (\mathbf{X}_U^T \mathbf{X}_U)^{-1} \mathbf{X}_U^T \mathbf{y}_U$$

- If $\lambda_N = o\left(\sqrt{N}\right)$, then LASSO shares this asymptotic equivalence
 - working model \neq true model, so “oracle properties” are not relevant
 - any advantages of LASSO are in finite samples

- Can also consider **adaptive lasso**, “A”: Zou (2006, *JASA*)
- Lasso is non-linear, so get weights either by
 - Wu and Sitter (2001) model calibration: “C”
 - Tibshirani (1996) ridge regression approximation to lasso coefficient estimates: “R”

nonlinear	calibration weights	ridge weights
LASSO	CLASSO	RLASSO
ALASSO	CALASSO	RALASSO

- Utah tree canopy cover data set
 - national pilot project conducted by FIA and the US Forest Service Remote Sensing Applications Center
 - $N = 4,151$ grid points on one Landsat scene in Utah
 - covers parts of 10 counties
- Response: $y_k =$ photo-interpreted tree canopy cover
 - relevant to forest management, fire modeling, air pollution mitigation, water temperature, and carbon storage
 - correlated with many other interesting responses
 - known for all $k \in U$ for this pilot study
 - very expensive to obtain!

- Auxiliary data x_k from the 2001 National Land Cover Database and Landsat-5 reflectance bands
 - transformed aspect, slope, topographic positional index, elevation, land cover and NLCD predicted tree canopy cover
 - about half are statistically significant in finite population regression of y on x
- Sampling designs:
 - simple random (SI) and stratified simple random (STSI) with counties as strata
 - samples sizes $n = 50$ and $n = 100$
 - equal allocation to counties \Rightarrow unequal probabilities
 - 2000 replicate samples from same fixed, finite population

- For all estimators, relative design bias was less than 2%
- Ratio of design MSE for each estimator to design MSE of full GREG estimator:

		SI		STSI	
		<i>n</i> = 50	<i>n</i> = 100	<i>n</i> = 50	<i>n</i> = 100
Model-Assisted:	LASSO	0.69	0.92	0.78	0.92
	ALASSO	0.76	0.94	0.80	0.94
	CLASSO	0.70	0.91	0.79	0.93
	CALASSO	0.77	0.94	0.83	0.94
	RLASSO	0.69	0.92	0.77	0.92
	RALASSO	0.71	0.97	0.69	0.95
Design Only:	HT	1.29	2.01	1.24	1.73

- Negative weights:

- calibration weights were negative in only 0.036% of all cases
- ridge regression weights in only 0.65%
- GREG weights varied from 1% to 14% negative weights

- Weight variation within samples and across samples:

Estimators	Weight Variances							
	$\overline{var}(\mathbf{w})$				$\overline{var}(w_j j \in s)$			
	SI		STSI		SI		STSI	
	$n = 50$	$n = 100$	$n = 50$	$n = 100$	$n = 50$	$n = 100$	$n = 50$	$n = 100$
CLASSO	158	18	5086	1216	155	18	837	92
CALASSO	170	18	5116	1216	166	18	894	93
RLASSO	496	47	6841	1322	441	46	3370	382
RALASSO	236	25	5059	1211	228	25	1059	125
GREG	5261	331	15349	1859	3810	298	12487	1310
HT	0	0	4816	1192	0	0	0	0

- Lasso survey regression estimators have useful potential
 - dominate MSE of GREG in small samples with large numbers of potential predictors
 - better weight properties: less variation and fewer negative weights
- Sophisticated computational methods can improve basic survey estimators
 - model-assisted framework gives straightforward recipe for incorporating complex methods
- Contact: jbreidt@stat.colostate.edu

Selected references on model-assisted survey regression estimation (by no means exhaustive, sorry!)

49

-
- Breidt, F. J., G. Claeskens, and J. D. Opsomer (2005). Model-assisted estimation for complex surveys using penalised splines. *Biometrika* 92, 831–846.
 - Breidt, F. J. and J. D. Opsomer (2000). Local polynomial regression estimators in survey sampling. *Annals of Statistics* 28, 1026–1053.
 - Breidt, F.J. and Opsomer, J.D. (2008). Endogenous post-stratification in surveys: classifying with a sample-fitted model. *Annals of Statistics* 36, 403–427.
 - Breidt, F. J. and J. D. Opsomer (2009). Nonparametric and semiparametric estimation in complex surveys. *Sample Surveys: Theory, Methods and Inference, Handbook of Statistics* 29, 103–119.
 - Breidt, F.J., Opsomer, J.D, Johnson, A.A. and Ranalli, M.G. (2007). Semiparametric model-assisted estimation for natural resource surveys. *Survey Methodology* 33, 35–44.
 - Cassel, C. M., C. E. Särndal, and J. H. Wretman (1976). Some results on generalized difference estimation and generalized regression estimation for finite populations. *Biometrika* 63, 615–620.
 - Dahlke, M., Breidt, F.J., Opsomer, J.D. and Van Keilegom, I. (2013). Nonparametric endogenous post-stratification estimation. *Statistica Sinica* 23, 189-211
 - Deville, J.C., and Särndal, C.E. (1992), Calibration Estimators in Survey Sampling, *Journal of the American Statistical Association* 87, 376-382.
 - Goga, C. (2005). Réduction de la variance dans les sondages en présence d'information auxiliaire: Une approche non paramétrique par splines de régression. *Canadian Journal of Statistics* 33, 163–180.

- McConville, K. S. and F. J. Breidt (2013). Survey design asymptotics for the model-assisted penalised spline regression estimator. *Journal of Nonparametric Statistics* (ahead-of-print), 1–19.
- Montanari, G. and M. Ranalli (2005). Nonparametric methods in survey sampling. *New Developments in Classification and Data Analysis 100*, 203–210.
- Montanari, G. E. and M. G. Ranalli (2005). Nonparametric model calibration estimation in survey sampling. *Journal of the American Statistical Association 100(472)*, 1429–1442.
- Montanari, G.E., and Ranalli, M.G.(2006), A Mixed Model-Assisted Regression Estimator that Uses Variables Employed at the Design Stage, *Statistical Methods and Applications 15*, 139–149.
- Montanari, G.E., and Ranalli, M.G. (2009). Multiple and ridge model calibration. In *Proc. Workshop on Calibration and Estimation in Surveys*. Statistics Canada.
- Opsomer, J. D., F. J. Breidt, G. G. Moisen, and G. Kauermann (2007). Model-assisted estimation of forest resources with generalized additive models (with discussion). *Journal of the American Statistical Association 102*, 400–416.
- Park, M. and Fuller, W.A. (2005), Towards nonnegative regression weights for survey samples. *Survey Methodology 31*, 85–93.
- Park, M. and Fuller, W.A. (2009), The mixed model for survey regression estimation, *Journal of Statistical Planning and Inference 139*, 1320–1331.
- Rao, J.N.K. and Singh, A.C. (1997). A ridge-shrinkage method for range-restricted weight calibration in survey sampling In *ASA Proc. Sect. Survey Res. Meth.*, American Statistical Association.
- Robinson, P. M. and C. E. Särndal (1983). Asymptotic properties of the generalized regression estimator in probability sampling. *Sankhya: The Indian Journal of Statistics, Series B 45*, 240–248.
- Särndal, C. E. (2007). The calibration approach in survey theory and practice. *Survey Methodology 33(2)*, 99–119.

- Särndal, C.-E., B. Swensson, and J. Wretman (1989). The weighted residual technique for estimating the variance of the general regression estimator of the finite population total. *Biometrika* 76, 527–537.
- Särndal, C.-E., B. Swensson, and J. Wretman (1992). *Model Assisted Survey Sampling*. New York: Springer-Verlag.
- Wang, L. and S. Wang (2011). Nonparametric additive model-assisted estimation for survey data. *Journal of Multivariate Analysis* 102, 1126–1140.
- Wu, C. and R. R. Sitter (2001). A model-calibration approach to using complete auxiliary information from survey data. *Journal of the American Statistical Association* 96, 185–193.