

B H F (1988) JASA

Prediction of county crop areas using survey and satellite data

- areas under corn and soybeans in 12 counties
- USDA Statistical Reporting Service
 - "exact" area of corn/soybeans in 37 segments in the 12 counties
 - Segment ≈ 250 hectares
- Satellite data available for all segments in the 12 counties
 - pixels are classified as corn/soybean
 - pixel ≈ 0.45 hectare
- Data

y_{ij} = number of hectares of corn/soybean in segment j of county i
 $i = 1, \dots, T = 12, j = 1, \dots, n_i$

x_{1ij} = number of pixels classified as corn in segment j in county i

x_{2ij} = - - - soybeans

- Model

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \gamma_i + e_{ij}$$

$$\gamma_i \sim \text{iid } N(0, \sigma_\gamma^2) \quad \begin{matrix} \text{random county} \\ \text{specific effect} \end{matrix}$$

$$\varepsilon_{ij} \sim \text{iid } N(0, \sigma_e^2)$$

- Matrix Formulation

$$Y = X\beta + Zv + \varepsilon$$

$$Y = (Y'_1, Y'_2, \dots, Y'_T)'$$

$$X = [1, z_1, z_2]$$

$$\beta = (\beta_0, \beta_1, \beta_2)'$$

$$Z = [z'_1, z'_2, \dots, z'_T]'$$

where z_i is $n_i \times 12$ and has rows equal to indicator vectors with 1 in position i

$$z = (z_1, \dots, z_{12})'$$

$$\text{var}(Y) = V = \text{blkdiag}(\sigma_e^2 I_{n_1} + \sigma_v^2 J_{n_1})$$

- Mean crop area per segment in county i conditional on γ_i is

$$\theta_i = \bar{x}_{i(P)}' \beta + \gamma_i$$

where $\bar{x}_{i(P)} = \frac{1}{N_i} \sum_{j=1}^{N_i} (1, x_{1ij}, x_{2ij})'$

and N_i = total number of segments in county i

$$\text{But } E(v_i | y_i) = g_i(\bar{y}_i - \tilde{x}'_i \beta)$$

$$\text{where } g_i = \frac{\sigma_v^2}{\sigma_v^2 + \sigma_e^2/\alpha_i}$$

BLUP for θ_i is

$$\tilde{\theta}_i = \tilde{x}'_{i(p)} \tilde{\beta} + \tilde{v}_i$$

$$\text{where } \tilde{\beta} = (x' V^{-1} x)^{-1} x' V^{-1} y$$

$$\text{and } \tilde{v}_i = g_i(\bar{y}_i - \tilde{x}'_i \tilde{\beta})$$

- Naive Prediction Variance

$$E\{(\tilde{\theta}_i - \theta_i)^2\} = \sigma_v^2(1-g_i) + c_i'(x' V^{-1} x)^{-1} c_i$$

$$\text{where } c_i = \tilde{x}'_{i(p)} - g_i \tilde{x}'_i$$

- In practice σ_v^2 and σ_e^2 also have to be estimated. The EBBLUE for θ_i is

$$\hat{\theta}_i = \tilde{x}'_{i(p)} \hat{\beta} + \hat{v}_i$$

- Survey Regression Estimate

→ fixed county specific effects.

$$\hat{\theta}_i = \bar{y}_i - (\tilde{x}'_i - \tilde{x}'_{i(p)})' \hat{\beta}$$

- BHF appendix - complicated formula for prediction variance that accounts for estimation of σ_v^2 and σ_e^2

Parametric Bootstrap Prediction Variances

- Simulate B "resamples" from the fitted model
 - $\gamma_i^* \sim N(0, \hat{\sigma}_\gamma^2)$
 - $e_{ij}^* \sim N(0, \hat{\sigma}_e^2)$
 - $y_{ij}^* = \bar{x}_{ij}' \hat{\beta} + \gamma_i^* + e_{ij}^*$
 - for each resample calculate predictions for the 12 countries; i.e.
- $$\hat{\theta}_{i,b}^* \quad i=1, \dots, 12, \quad b=1, \dots, B$$
- Notice that the true values for each resample are known
- $$\theta_{i,b}^* = \bar{x}_{i(p)}' \hat{\beta} + \gamma_i^*$$
- $$\hat{\theta}_{i,b}^* = \bar{x}_{i(p)}' \hat{\beta}^* + \hat{\gamma}_i^*$$
- Bootstrap prediction variance for county i is

$$\frac{1}{B} \sum_{b=1}^B (\hat{\theta}_{i,b}^* - \theta_{i,b}^*)^2$$

The BLUP and Penalized Least Squares

$$Y = X\beta + Z\alpha + E$$

$$u \sim N(0, D), \quad e \sim N(0, \sigma_e^2 I)$$

- (γ, u) jointly normal with log density

$$-\frac{1}{2} \log |2\pi \sigma_e^2 I| - \frac{1}{2} \log (2\pi D)$$

$$-\frac{1}{2} \left[\frac{1}{\sigma_e^2} \| \gamma - X\beta - Zu \|^2 + u'D^{-1}u \right]$$

- Maximizing jointly with respect to β and u is equivalent to minimizing

$$\frac{1}{\sigma_e^2} \| \gamma - X\beta - Zu \|^2 + u'D^{-1}u$$

This results in

$$\hat{\beta} = (X'\Sigma^{-1}X)^{-1} X'\Sigma^{-1}\gamma \quad \text{BLUE}$$

$$\hat{u} = C'\Sigma^{-1}(\gamma - X\hat{\beta}) \quad \text{BLUP}$$

where $\Sigma = \sigma_e^2 I + ZDZ'$ and $C = ZD$

(Henderson, 1950)

Penalized Spline Regression

(Reference: Ruppert, Wand, Carroll, 2003
"Semi parametric Regression")

- Want to fit a model of the form

$$y_i = f(x_i) + e_i$$

to data $(x_i, y_i) \quad i=1, \dots, n$

- linear regression

$$f(x) = \beta_0 + \beta_1 x$$

- least squares estimation - minimize

$$\|y - X\beta\|^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

- linear spline model

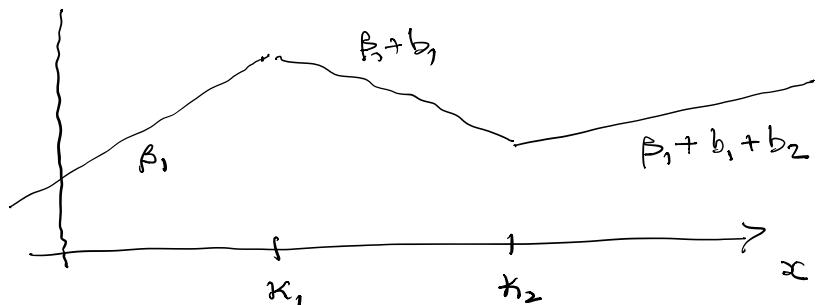
$$f(x) = \beta_0 + \beta_1 x + \sum_{j=1}^K b_j (x - x_j)_+$$

where $\min(x_i) < x_1 < x_2 < \dots < x_K < \max(x_i)$

are "knots" and

$$(x - x_j)_+ = \max(0, x - x_j)$$

e.g. $f(x) = \beta_0 + \beta_1 x + b_1 (x - x_1)_+ + b_2 (x - x_2)_+$



least squares fit with

$$X = (1, x, (x - 1x_1)_+, \dots, (x - 1x_K)_+)$$

$f(x)$ has erratic changes in slope

- Penalized Least Squares

minimize $\|Y - X\beta\|^2$ subject to $\beta' A \beta \leq C$

for some psd A and $C > 0$. This is equivalent to minimizing

$$\|Y - X\beta\|^2 + \lambda \underbrace{\beta' A \beta}_{\text{penalty term}}$$

for some $\lambda > 0$

Taking derivatives

$$\begin{aligned} \frac{\partial}{\partial \beta} & \left[(Y - X\beta)' (Y - X\beta) + \lambda \beta' A \beta \right] \\ &= -2X'(Y - X\beta) + 2\lambda A\beta \end{aligned}$$

The solution is

$$\hat{\beta} = (X'X + \lambda A)^{-1} X' Y$$