

SRS
1981

1981

PREDICTION OF COUNTY CROP AREAS USING SURVEY AND SATELLITE DATA

George E. Battese, University of New England
Wayne A. Fuller, Iowa State University

INTRODUCTION

The estimation of parameters for small areas has received considerable attention in recent years. A comprehensive review of research in small-area estimation is given by Purcell and Kish (1979). Agencies of the Federal Government have been significantly involved in this research to obtain estimates of such items as population counts, unemployment rates, per capita income, health needs, etc., for states and local government areas. Acts of the U.S. Congress (e.g., Local Fiscal Assistance Act of 1972 and the National Health Planning and Resources Development Act of 1974) have created a need for accurate small-area estimates. Research in this area is illustrated in such papers as DiGaetano et al. (1980), Fay and Herriot (1979), Erickson (1974), Gonzalez (1973) and Gonzalez and Hoza (1978). Fay and Herriot (1979) outline the approach used by the U.S. Bureau of the Census which is based upon James-Stein estimators [see Efron and Morris (1973), James and Stein (1961)].

Recently the U.S. Department of Agriculture (U.S.D.A.) has been investigating the use of LANDSAT satellite data to improve its estimates of crop areas for Crop Reporting Districts and to develop estimates for individual counties. The methodology used in some of these studies is presented in Cárdenas, Blanchard and Craig (1978), Hanuschak et al. (1979), and Sigman et al. (1978). In these studies ground observations from the U.S.D.A.'s June Enumerative Survey for sample segments are regressed on the corresponding satellite data for given strata. County estimates obtained by the regression approach generally have smaller estimated variances than those for the traditional "direct expansion approach" using survey data only.

In this paper we consider the prediction of crop areas in counties for which survey and satellite data are available. It is assumed that for sample counties, reported crop areas are obtained for a sample of area segments by interviewing farm operators. We assume that data for more than one sample segment are available for several sample counties. In addition, we assume that for each sample segment and county, satellite data are obtained and the crop cover classified for each pixel. A pixel (an acronym for "picture element") is the unit for which satellite information is recorded and is about 0.45 hectares in area. Predictors for county crop areas are obtained under the assumption that the nested-error regression model defines the relationship between the survey and satellite data.

NESTED-ERROR MODEL AND PREDICTORS

Consider the model

$$Y_{ij} = x_{ij}\beta + u_{ij}, \quad i=1, \dots, t; \quad j=1, \dots, n_i; \quad (1)$$

and

$$u_{ij} = v_i + e_{ij}, \quad (2)$$

where Y_{ij} is the reported area in the given crop for the j -th sample segment of the i -th county as recorded in the sample survey involved; n_i is the number of sample segments observed in the i -th sample county; x_{ij} is a $(1 \times k)$ vector of values of explanatory variables which are functions of the satellite data; and β is a $(k \times 1)$ vector of unknown parameters. The random errors, $v_i, i = 1, 2, \dots, t$, are assumed to be N.I.D. $(0, \sigma_v^2)$ independent of the e_{ij} 's, which are assumed to be N.I.D. $(0, \sigma_e^2)$. From these assumptions it follows that the covariance structure of the errors of the model is given by

$$E(u_{ij}u_{i'j'}) = \begin{cases} \sigma_v^2 + \sigma_e^2, & \text{if } i=i' \text{ and } j=j' \\ \sigma_v^2, & \text{if } i=i' \text{ and } j \neq j' \\ 0, & \text{if } i \neq i' \end{cases} \quad (3)$$

This model specifies that the reported crop areas for segments within a given county are correlated, and that the covariances are the same for all counties, but that the reported crop areas for different counties are not correlated. Efficient estimation of the nested-error model is discussed in Fuller and Battese (1973).

We consider that for each sample county, the mean crop area per segment is to be predicted. These are conditional means that are denoted by $\mu_i, i = 1, 2, \dots, t$, where

$$\mu_i = \bar{x}_{i(p)}\beta + v_i, \quad (4)$$

where $\bar{x}_{i(p)} \equiv N_i^{-1} \sum_{j=1}^{N_i} x_{ij}$, the mean of the x_{ij} 's for the N_i population segments in the i -th county, is assumed known. Note that μ_i is the conditional mean of Y_{ij} for the i -th sample county, given the population mean of the x_{ij} -values.

It is noted that the above problem is a special case of the estimation of a linear combination of fixed effects and realized values of random effects [see Harville (1976), (1979), and Henderson (1975)]. Although some of our results are obtained as special cases of the general

1981 ASA Survey Section Proceedings

results are given in these papers, we derive the predictors involved by first considering that the elements of β (as well as the variance components, $\sigma_v^2 > 0$ and $\sigma_e^2 > 0$ are known. Considering the prediction of the county effects, v_i , $i = 1, 2, \dots, t$, is motivation for the predictors of the county means to be presented later.

(a) Prediction When β is Known

If the parameters of the model (1) are known, then the random errors, u_{ij} , are observable. The sample mean of the random errors for the i -th county, $\bar{u}_i \equiv n_i^{-1} \sum_{j=1}^{n_i} u_{ij} = v_i + \bar{e}_i$. has unconditional mean 0 and variance $\sigma_v^2 + n_i^{-1} \sigma_e^2$. The conditional mean, $E(\bar{u}_i | v_i)$, is equal to v_i and the conditional variance, $V(\bar{u}_i | v_i)$, is equal to $n_i^{-1} \sigma_e^2$. Thus, the sample mean, \bar{u}_i , is a conditionally unbiased predictor for v_i , $i=1, 2, \dots, t$.

We consider the class of linear predictors for v_i that is defined by

$$\hat{v}_i^{(\delta)} \equiv \delta_i \bar{u}_i,$$

where δ_i is a constant such that $0 \leq \delta_i \leq 1$. The error in this predictor is given by

$$\hat{v}_i^{(\delta)} - v_i = -(1-\delta_i)v_i + \delta_i \bar{e}_i. \quad (5)$$

and so the mean squared error of the predictor is

$$E[\hat{v}_i^{(\delta)} - v_i]^2 = (1-\delta_i)^2 \sigma_v^2 + \delta_i^2 n_i^{-1} \sigma_e^2. \quad (6)$$

It is easily verified that (6) is equal to

$$E[\hat{v}_i^{(\delta)} - v_i]^2 = (\sigma_v^2 + n_i^{-1} \sigma_e^2) (\delta_i - \gamma_i)^2 + (1-\gamma_i) \sigma_v^2,$$

where γ_i is defined by

$$\gamma_i = \sigma_v^2 (\sigma_v^2 + n_i^{-1} \sigma_e^2)^{-1}. \quad (7)$$

Thus, the best linear predictor of v_i is

$$\hat{v}_i^{(\gamma)} \equiv \gamma_i \bar{u}_i \text{ and its mean squared error is}$$

$$E[\hat{v}_i^{(\gamma)} - v_i]^2 = (1-\gamma_i) \sigma_v^2 \equiv \gamma_i n_i^{-1} \sigma_e^2. \quad (8)$$

However, under the assumption of normality of v_i and e_{ij} , it follows that $E(v_i | \bar{u}_i) = \gamma_i \bar{u}_i$ and so $\hat{v}_i^{(\gamma)}$ is the predictor with minimum mean squared error.

Since the expectation of the prediction error (5), conditional on v_i , is $-(1-\delta_i)v_i$, then the mean of the squared conditional bias (hereafter referred to as the mean squared bias) of the predictor is

$$E[E(v_i^{(\delta)} | v_i) - v_i]^2 = (1-\delta_i)^2 \sigma_v^2. \quad (9)$$

We consider that the mean squared bias (9) may be of basic importance in the consideration of predictors for v_i . In fact, we assume that information is available such that the mean squared bias of predictors is required to be no larger than a predetermined constant. The best predictor with constrained mean squared bias is stated in Theorem 1.

Theorem 1. If model (1)-(2) holds and the parameters are known, then, in the class of linear predictors, $\hat{v}_i^{(\delta)} \equiv \delta_i \bar{u}_i$, for which the mean squared bias is constrained by Δ_i , the best predictor is defined by

$$\hat{v}_i^{(\gamma^*)} = \begin{cases} \gamma_i \bar{u}_i, & \text{if } \Delta_i \geq (1-\gamma_i)^2 \sigma_v^2 \\ \gamma_i^* \bar{u}_i, & \text{if } \Delta_i < (1-\gamma_i)^2 \sigma_v^2 \end{cases} \quad (10)$$

where $\gamma_i^* = 1 - (\Delta_i / \sigma_v^2)^{1/2}$; and therefore predictors satisfying $0 \leq \delta_i < \gamma_i$ are inadmissible.

(b) Prediction When β is Unknown

Returning to the prediction of the conditional county means, $\mu_i \equiv \bar{x}_i(p)\beta + v_i$, $i = 1, 2, \dots, t$, it is seen that the problem is that of predicting the sum of v_i and a linear function of unknown parameters. We consider the class of predictors defined by

$$\hat{\mu}_i^{(\delta)} \equiv \bar{x}_i(p)\hat{\beta} + \delta_i (\bar{y}_i - \bar{x}_i\hat{\beta}), \quad (11)$$

where $\hat{\beta}$ is the best linear unbiased estimator for β (assuming again that $\sigma_v^2 > 0$ and $\sigma_e^2 > 0$ are known); $\bar{y}_i \equiv n_i^{-1} \sum_{j=1}^{n_i} y_{ij}$; $\bar{x}_i \equiv n_i^{-1} \sum_{j=1}^{n_i} x_{ij}$; and δ_i is a constant such that $0 \leq \delta_i \leq 1$.

For $\delta_i = 0$, the predictor (11) is

$$\hat{\mu}_i^{(0)} \equiv \bar{x}_i(p)\hat{\beta},$$

which is referred to as the "regression predictor." For $\delta_1 = 1$, the predictor is

$$\begin{aligned} \hat{\mu}_1^{(1)} &\equiv \bar{x}_{1(p)} \hat{\beta} + (\bar{y}_1 - \bar{x}_{1.} \hat{\beta}) \\ &= \bar{y}_1 + (\bar{x}_{1(p)} - \bar{x}_{1.}) \hat{\beta}, \end{aligned}$$

which is referred to as the "adjusted survey predictor." This predictor adjusts the survey sample mean, \bar{y}_1 , to account for the sample mean of the regressors, $\bar{x}_{1.}$, differing from the population mean, $\bar{x}_{1(p)}$.

The error in the predictor (11) is expressed by

$$\begin{aligned} \hat{\mu}_1^{(\delta)} - \mu_1 &= [-(1-\delta_1)v_1 + \delta_1 e_{1.}] \\ &\quad + (\bar{x}_{1(p)} - \delta_1 \bar{x}_{1.}) (\hat{\beta} - \beta), \end{aligned} \quad (12)$$

where the first term is the prediction error (5) for the case when β is known and the second term arises in the estimation of β . The mean squared error for the general predictor (11) and the best linear predictor are stated in the next theorem.

Theorem 2. If model (1)-(2) holds, where σ_v^2 and σ_e^2 are known positive constants, then the mean squared error for the predictor $\hat{\mu}_1^{(\delta)}$ is

$$\begin{aligned} E[\hat{\mu}_1^{(\delta)} - \mu_1]^2 &= [(1-\delta_1)^2 \sigma_v^2 + \delta_1^2 n_1^{-1} \sigma_e^2] \\ &\quad + 2(\delta_1 - \gamma_1) (\bar{x}_{1(p)} - \delta_1 \bar{x}_{1.}) V(\hat{\beta}) \bar{x}_{1.}' \\ &\quad + (\bar{x}_{1(p)} - \delta_1 \bar{x}_{1.}) V(\hat{\beta}) (\bar{x}_{1(p)} - \delta_1 \bar{x}_{1.})', \end{aligned} \quad (13)$$

where $V(\hat{\beta})$ is the covariance matrix for $\hat{\beta}$. Furthermore, the mean squared error is a minimum when $\delta_1 = \gamma_1 \equiv \sigma_v^2 (\sigma_v^2 + n_1^{-1} \sigma_e^2)^{-1}$.

It can be shown that the expectation of the prediction error (12), conditional on the realized random effects, $v = (v_1, v_2, \dots, v_t)'$ is

$$\begin{aligned} E[\hat{\mu}_1^{(\delta)} | v] - \mu_1 &= -(1-\delta_1)v_1 \\ &\quad + (\bar{x}_{1(p)} - \delta_1 \bar{x}_{1.}) V(\hat{\beta}) \sum_{j=1}^t \bar{x}_{j.}' v_j \gamma_j / \sigma_v^2. \end{aligned} \quad (14)$$

From this result the mean squared bias of $\hat{\mu}_1^{(\delta)}$,

denoted by $MSB(\hat{\mu}_1^{(\delta)})$, and the best constrained predictor are obtained, as stated in the following theorem.

Theorem 3. If model (1)-(2) holds, where σ_v^2 and σ_e^2 are known positive constants, then the mean squared bias of the predictor $\hat{\mu}_1^{(\delta)}$ is

$$\begin{aligned} E[E(\hat{\mu}_1^{(\delta)} | v) - \mu_1]^2 &= (1-\delta_1)^2 \sigma_v^2 \\ &\quad - 2(1-\delta_1) \gamma_1 (\bar{x}_{1(p)} - \delta_1 \bar{x}_{1.}) V(\hat{\beta}) \bar{x}_{1.}' \\ &\quad + \sum_{j=1}^t [(\bar{x}_{1(p)} - \delta_1 \bar{x}_{1.}) V(\hat{\beta}) \bar{x}_{j.}' \gamma_j]^2 / \sigma_v^2. \end{aligned} \quad (15)$$

Furthermore, if the mean squared bias is constrained by Δ_1 , then the best constrained predictor is defined by

$$\begin{aligned} \hat{\mu}_1^{(\gamma^*)} &= \bar{x}_{1(p)} \hat{\beta} + \gamma_1 (\bar{y}_1 - \bar{x}_{1.} \hat{\beta}), \text{ if } \Delta_1 > MSB(\hat{\mu}_1^{(\gamma)}) \\ &= \bar{x}_{1(p)} \hat{\beta} + \gamma_1^* (\bar{y}_1 - \bar{x}_{1.} \hat{\beta}), \text{ if } \Delta_1 < MSB(\hat{\mu}_1^{(\gamma)}) \end{aligned}$$

where γ_1^* is the root of $\Delta_1 = E[E(\hat{\mu}_1^{(\delta)} | v) - \mu_1]^2$ with smaller mean squared error.

The mean squared bias (15) has positive derivative with respect to δ_1 and is generally expected to be monotone decreasing as δ_1 increases from zero to one.

ESTIMATION OF VARIANCES

When the variance components, σ_v^2 and σ_e^2 , are unknown, different estimators can be used. Harville (1977) contains a discussion of estimation methods for component-of-variance models. We use the fitting-of-constants estimators presented in Fuller and Battese (1973) for the nested-error model (1)-(2). By use of normal theory, the variances of these variance components are obtained and presented in Battese and Fuller (1981).

These estimators for the variances and covariances of the variance estimators are necessary for inference about σ_v^2 and σ_e^2 and for obtaining approximate generalized least-squares estimates for the variances when prior information is available. The county predictors defined above are approximated by replacing the variances, σ_v^2 and σ_e^2 , with their corresponding sample estimates.

If all n_i are the same, it is possible to use the results of Efron and Morris (1973) to show that the estimation of σ_v^2 will increase the average squared error by an amount approximately equal to

$$2(t-1)^{-1}(\sigma_v^2 + n_1^{-1}\sigma_e^2)^{-1} n_1^{-2} \sigma_e^4.$$

Estimating σ_e^2 adds a term to the variance that is approximately equal to

$$\left[\frac{n_1^{-1}(\sigma_v^2 + t_0^{-1}\sigma_e^2)}{(\sigma_v^2 + n_1^{-1}\sigma_e^2)^2} \right]^2 W(\hat{\sigma}_e^2) W(\bar{u}_{1.}),$$

where $t_0 = (t-1)^{-1} [n - n^{-1} \sum_{i=1}^t n_i^2]$.

EMPIRICAL RESULTS

We consider prediction of areas of soybeans for 12 counties in North-Central Iowa, based on data for 1978. The Economics and Statistics Service of the U.S. Department of Agriculture determined the area of soybeans in 37 area sampling units (segments) in the 12 counties during the June Enumerative Survey in 1978. The segments are about 259 hectares (or one square mile) in area. The numbers of pixels classified as soybeans in these area segments were determined from the NASA's LANDSAT satellites during passes over Iowa in August and September 1978.

To obtain predictions of crop areas we assume the simple model,

$$Y_{ij} = \beta_0 + \beta_1 X_{ij} + u_{ij}, \quad i=1,2,\dots,12; \quad j=1,2,\dots,n_i.$$

where Y_{ij} is the number of hectares of soybeans in the j -th sample segment of the i -th county as recorded in the June Enumerative Survey in 1978; X_{ij} is the number of pixels of soybeans for the j -th sample segment of county i . The parameter estimates are obtained by use of the nested-error software of SUPER CARP [Hidioglou, Fuller, and Hickman (1980)]. The parameter estimates and their estimated standard errors (in parentheses) are:

$$\hat{Y}_{ij} = -3.8 + 0.475 X_{ij}, \quad \text{where}$$

(9.3) (0.040)

$$\hat{\sigma}_v^2 = 250 \quad \text{and} \quad \hat{\sigma}_e^2 = 184.$$

(142) (53)

The estimate for the intercept parameter is not significantly different from zero. The among-county variance estimate, σ_v^2 , is significant at the 5% level. The fact that σ_v^2 and

σ_e^2 were estimated was ignored in computing the standard errors of the remaining parameters.

Given the preceding results, the predictions for the mean hectares of soybeans per segment in the several counties are listed in Table 1 for the different predictors discussed in preceding sections. Also presented are the sample mean of the reported soybean hectares for the June Enumerative Survey. The square root of the estimated mean squared error is given in parentheses below the corresponding prediction. The population and sample county means for the number of pixels classified as soybeans are presented in Table 2, together with the population number of segments.

It is evident from Table 1 that the use of the satellite data to obtain predictors is much more efficient than using only the reported crop areas from the June Enumerative Survey. The mean squared errors of the sample mean of the reported hectares are relatively large for the individual counties. For predicting soybean areas, the regression predictor is always less efficient than the adjusted survey predictor because the variance among counties is the dominant term in the total variance.

The best predictor, $\hat{\mu}_1^{(\gamma)}$, has mean squared error that is considerably smaller than that for the regression predictor, $\hat{\mu}_1^{(0)}$, especially when several sample segments are available in a county. The ratio of the mean squared error for the best predictor to that for the regression predictor is a rather complicated function of the variances, the values of the x -variables, and the sample sizes. However, for the soybean data, the values of this ratio for the different counties varied little for particular numbers of sample segments. The square root of the average value of the ratio is presented in Table 3 for the different values of the sample sizes. Although these statistics are based on different numbers of observations and should be interpreted with caution, they show an interesting pattern. As the number of sample segments increases, the relative root mean squared error decreases, but at a declining rate. This is due to the fact that the mean squared error for the best predictor decreases markedly as the number of sample segments increases, but that for the regression predictor does not. The decreases in the relative root mean squared error with increasing numbers of sample segments are quite substantial for soybeans. Furthermore, these data would suggest that obtaining data for a few sample segments in more counties is likely to result in greater precision of prediction than obtaining more data for fewer counties.

CONCLUSIONS

The nested-error regression model with satellite data as the auxiliary variable offers a promising approach to prediction of crop areas

Table 1: Predicted Hectares of Soybeans per Segment for Twelve Iowa Counties

County i	γ_i	Predictions			
		$\hat{\mu}_i^{(0)}$	$\hat{\mu}_i^{(\gamma)}$	$\hat{\mu}_i^{(1)}$	\bar{y}_i
Cerro Gordo	0.58	86.4 (15.6)	78.2 (11.0)	72.1 (13.7)	8.1 (31.4)
Franklin	0.80	85.6 (15.3)	66.1 (7.1)	61.4 (7.8)	52.5 (18.2)
Hamilton	0.58	89.7 (15.7)	93.3 (10.5)	95.9 (13.6)	106.0 (31.4)
Hancock	0.87	90.7 (15.2)	100.5 (5.8)	101.9 (6.2)	117.5 (14.1)
Hardin	0.89	80.4 (15.2)	74.4 (5.4)	73.7 (5.7)	89.8 (12.8)
Humboldt	0.73	100.9 (15.6)	81.8 (8.7)	74.7 (9.9)	35.1 (22.2)
Kossuth	0.87	93.5 (15.2)	119.3 (5.7)	123.1 (6.1)	117.8 (14.1)
Pocahontas	0.80	113.7 (15.2)	113.2 (7.1)	113.1 (7.8)	118.7 (18.2)
Webster	0.84	113.7 (15.1)	109.9 (6.3)	109.2 (6.8)	113.0 (15.7)
Winnebago	0.80	84.3 (15.3)	97.6 (7.1)	100.8 (7.9)	88.6 (18.2)
Worth	0.58	93.8 (15.7)	87.2 (10.6)	82.3 (13.6)	103.6 (31.4)
Wright	0.80	101.5 (15.3)	112.8 (7.2)	115.6 (8.0)	97.8 (18.2)

Table 2: Pixel Data for Soybeans in Twelve Iowa Counties

County	No. of segments in pop'n sample		Pop'n mean	Sample mean
Cerro Gordo	545	1	189.70	55.00
Franklin	564	3	188.06	169.33
Hamilton	566	1	196.65	218.00
Hancock	569	5	198.66	231.40
Hardin	556	6	177.05	210.83
Humboldt	424	2	220.22	137.00
Kossuth	965	5	204.61	193.60
Pocahontas	570	3	247.13	259.00
Webster	687	4	247.09	255.00
Winnebago	402	3	185.37	159.67
Worth	394	1	205.28	250.00
Wright	567	3	221.36	184.00

Table 3: Averages of the Ratio of the Root Mean Squared Error for the Best Predictor to that for the Regression Predictor

Number of sample segments, n_i	$\{MSE(\hat{\mu}_i^{(\gamma)})/MSE(\hat{\mu}_i^{(0)})\}^{1/2}$
1	0.68
2	0.56
3	0.47
4	0.42
5	0.38
6	0.38

for counties. A reasonably large number of counties is required for the satisfactory estimation of the among-county variance. The U.S. Department of Agriculture plans to implement the software of the nested-error approach for the prediction of county crop areas in the next crop year.

REFERENCES

- Battese, G. E., and Fuller, W. A. (1981), "Prediction of County Crop Areas Using Survey and Satellite Data," unpublished paper, Statistical Laboratory, Iowa State University, Ames.
- Cárdenas, M., Blanchard, M. M., and Craig, M. E. (1978), On the Development of Small Area Estimators Using LANDSAT Data as Auxiliary Information, Economics, Statistics, and Cooperatives Service, U.S. Department of Agriculture, Washington, D.C.
- DiGaetano, R., MacKenzie, E., Waksberg, J., and Yaffe, R. (1980), "Synthetic Estimates for Local Areas from the Health Interview Survey," 1980 Proceedings of the Section on the Survey Research Methods, American Statistical Association.
- Efron, B., and Morris, C. (1973), "Stein's Estimation Rule and Its Competitors-An Empirical Bayes Approach," Journal of the American Statistical Association, 68, 117-130.
- Ericksen, E. P. (1974), A Regression Method for Estimating Population Changes of Local Areas. Journal of the American Statistical Association, 69, 867-875.
- Fay, R. E., and Herriot, R. (1979), "Estimates of Income for Small Places: An Application of James-Stein Procedures to Census Data," Journal of the American Statistical Association, 74, 269-277.
- Fuller, W. A., and Battese, G. E. (1973), "Transformations for Estimation of Linear Models with Nested-Error Structure," Jour-

nal of the American Statistical Association,
68, 626-632.

Gonzalez, M. E. (1973), "Use and Evaluation of Synthetic Estimates," 1973 Proceedings of the Social Statistics Section, American Statistical Association, 33-36.

Gonzalez, M. E., and Hoza, C. (1978), "Small Area Estimation with Application to Unemployment and Housing Estimates," Journal of the American Statistical Association, 73, 7-15.

Hanuschak, G., Sigman, R., Craig, M., Ozga, M., Luebbe, R., Cook, P., Kleweno, D., and Miller, C. (1979), Obtaining Timely Crop Area Estimates Using Ground-Gathered and LANDSAT Data, Technical Bulletin No. 1609, Economics, Statistics, and Cooperatives Service, U.S. Department of Agriculture, Washington, D.C.

Harville, D. A. (1976), "Extension of the Gauss-Markov Theorem to Include the Estimation of Random Effects," The Annals of Statistics, 4, 384-395.

Harville, D. A. (1977), "Maximum Likelihood Approaches to Variance Component Estimation and Related Problems," Journal of the American Statistical Association, 72, 320-338.

Harville, D. A. (1979), "Some Useful Representations for Constrained Mixed-Model Estimation," Journal of the American Statistical Association, 74, 200-206.

Henderson, C. R. (1975), "Best Linear Unbiased

Estimation and Prediction Under a Selection Model," Biometrics, 31, 423-447.

Hidiroglou, M. A., Fuller, W. A., and Hickman, R. D. (1980), SUPER CARP, Sixth Edition, Survey Section, Statistical Laboratory, Iowa State University, Ames, Iowa.

James, W., and Stein, Charles (1961), "Estimation with Quadratic Loss," Proceedings of the Fourth Berkeley Symposium of Mathematical Statistics and Probability, Vol. 1, University of California Press, 361-379.

Purcell, N. J., and Kish, L. (1979), "Estimation for Small Domains," Biometrics, 35, 365-384.

Rao, C. R. (1965), Linear Statistical Inference and its Applications, Wiley, New York.

Sigman, R. S., Hanuschak, G. A., Craig, M. E., Cook, P. W., and Cárdenas, M. (1978), "The Use of Regression Estimation with LANDSAT and Probability Ground Sample Data," 1978 Proceedings of the Section on Survey Research Methods, American Statistical Association, 165-168.

ACKNOWLEDGMENTS

This research was partly supported by Research Agreement 58-319T-1-0054X with the Economics and Statistics Service of the U.S. Department of Agriculture. We thank Cheryl Auer for writing the computer programs for the empirical analyses and Rachel Harter for helpful discussions on earlier drafts of the paper.

1981

DISCUSSION

Raj S. Chhikara, Old Dominion University

Two of the three papers related to remote sensing, one by Professors Battese and Fuller and another by Amis, et al, deal with the problem of predicting county crop acreages using the regular USDA survey data in conjunction with satellite (Landsat) data. Currently in USDA, the regression estimator obtained by regressing the survey reported crop acreages onto those determined from Landsat data for sample segments in the area of interest is used for estimation of crop acreages. Each of these two papers discusses methods by which county crop acreages can be predicted more precisely with the utilization of satellite data in addition to the regular survey data. Professors Battese and Fuller have proposed a predictor which has the desirable property of minimum mean square error; whereas Amis, et al, have investigated empirically methods of improving classification of Landsat data.

First I discuss the paper by Battese and Fuller. It generalizes the classical regression estimator and gives a class of linear predictors assuming a nested-error regression model. This is achieved by weighting the classical estimator where weights depend upon the sample size, the county and error variance components. The minimum mean square error predictor is derived under the assumption of known variances. Consideration is given to the problem of bias and a modified predictor is suggested when the mean squared bias is desired not to exceed a threshold value. The method is applied to predict corn and soybean acreages for 12 counties in North-Central Iowa. The best predictor, standard regression predictor and two other predictors are computed. Numerical results show that the estimated mean square error is minimum for the best predictor.

This is an excellent paper showing how the present USDA crop acreage estimates can be improved upon at the county level. This is of course to meeting certain underlying assumptions. I have a few concerns regarding the assumptions made in the paper.

First, which is of minor importance, is the assumption of known variances. In general, variances are unknown and only their estimates can be available for use in prediction. So the optimal property of their predictor holds only conditionally.

My major concern is the assumption of known Landsat crop pixels in a segment. These crop pixels are estimated, sometimes with gross errors, and hence, the auxiliary variable is subject to measurement error. In the present context of classification of Landsat data with many limitations in adequately training a classifier, the measurement error is not necessarily uniform over an entire Landsat scene, and thus, it may introduce bias in the predictor. If this bias is considerable and dominates other errors, it should be investigated and, if necessary, a predictor which is at least approximately unbiased be constructed.

Lastly, the paper gives a predictor for the overall mean crop acreage for a collection of counties. This predictor is not necessarily the optimum one despite of its being the linear combination of the best predictors for individual

counties. The investigators may consider constructing a predictor which has the minimum mean square error. The adjustment proposed in the paper seems artificial and does not necessarily improve the overall mean predictor except matching the large area crop acreage to the aggregated county crop acreages.

The paper by Amis, et al, focuses on the problem of classification of Landsat data to estimate the number of pixels for different crops and the extent to which the errors in classification affect the crop acreage estimation error. The classification system presently used--called EDITOR--may be iterated several times in selecting data for training of a classifier, if necessary, to achieve maximum value for the square of correlation coefficient, r^2 , for the sample segments for which both survey reported crop acreages and Landsat data are available. The precision of a crop acreage estimate is based on the variance estimated by the residual mean square error times $(1-r^2)$. Since the value of r^2 for the sample segments which are used in training the classifier is expected to be higher than that for the segments which are not used in such training, a smaller value of r^2 is expected for the entire scene and hence, the gain in precision is likely to be overestimated.

The empirical study is conducted to evaluate overestimation of the gain in precision using 33 segments available from Missouri. Alternative clustering and classification techniques are also considered to seek improvement in the classification performance. The set of segments are treated in three different ways: (1) All 33 segments are used in training the classifier and in obtaining the regression equation, (2) 25 segments for training the classifier and 8 segments for an independent test set for the classification and regression, (3) jackknifing with 30 segments for training the classifier and 3 segments for the test repeated 11 times. Based on the test results it was concluded that:

- (i) The classification error rates were higher for the test segments as compared to those for the segments used in training of the classifier.
- (ii) The value of r^2 were smaller for the test segments as compared to those for the segments used in training of the classifier, implying that the gain in precision of a crop acreage estimate is overestimated.
- (iii) Use of the alternative clustering method--called CLASSY--resulted in a smaller mean square error as well as in lower classification rates. The use of CLASSY was recommended to improve upon the present method of crop acreage estimation.

Most of the statistical analysis was based on comparative tests using the Hotelling's T^2 statistic. No specific statistical inference was made on the overestimation of precision and on the determination of bias in a crop acreage estimate resulting from the classification errors. Since each of these issues is equally important in the evaluation of the present approach to crop acreage

1981 ASA ~~Survey~~ Survey Section Proceedings

estimation, different analyses of data should be constructed. Because of a limited data set these analyses may not be conclusive, yet the available data can be used to plan another empirical study.

It seems the present study is very limited and needs to be extended to a larger region with a larger set of test segments. It will be beneficial to include in any new empirical study the testing of the linear predictor proposed by Battese and Fuller to investigate how well it performs in the presence of classification errors.

The paper by Ron Fesco addresses the practical problem of stratification and area frame develop-

ment for a large-scale crop survey using satellite data. He has presented a procedure for obtaining a land-use stratification using information derived mostly from the Landsat data. An automated area sampling frame, given in a digitized form and which is flexible enough to permit changes in sampling unit size, is proposed and discussed. I did not get a chance to study his gridding system in details. I hope that the stratification and area frame resulting from the proposed method is tested for its efficiency for a region in the U.S. before its implementation on a large-scale.