

Alexander Buyantuyev and Jianguo Wu  
 School of Life Sciences, Arizona State University, Tempe

## ABSTRACT

We conducted error analysis on the statistical relationships between projected foliage cover obtained for the CAPLTER 200-point survey and vegetation cover estimated from remote sensing data.

We developed alternative bivariate linear regressions and used bootstrap and jackknife resampling to test these models and compare uncertainties (variability of the slope term). Although these regressions showed similar ability to explain variation in data, they had noticeably different slopes, reflecting the differences in the fundamental assumptions of these models. The models were then compared using estimated overall bias, root mean square errors (RMSE), standard error (SE), and variance ratios. RMA regression provided a better fit for desert sites, but standard OLS was superior for urban plots. However, the use of traditional OLS, which assumes no measurement errors in predictor variables, is considered flawed in principle. Agricultural vegetation was modeled with RMA regression.

Accuracy analysis was conducted by comparing vegetation cover observed on high spatial resolution color aerial photography with cover predicted from best statistical models. A total of 175 validation sites were stratified by three major land uses and randomly placed within urban, desert, and agricultural land uses. We examined correlations between predicted and observed cover and used RMSE and SE to assess the overall accuracy of vegetation cover prediction. The results suggest good agreement for urban land use, somewhat less accurate predictions for agricultural land use, and acceptable overall accuracy for desert vegetation.

## DEVELOPMENT OF MAP OF VEGETATION COVER

Three Landsat ETM+ images acquired at three dates (3/18/00, 4/19/00, and 5/21/00) were atmospherically corrected, spatially re-registered to high-resolution air photo, and used to derive vegetation indices (NDVI and SAVI) and subpixel fractions of vegetation endmember (UNMIX). NDVI, SAVI, and UNMIX image pixels corresponding to 200 survey plots were sampled and subdivided into three groups based on CAPLTER land use classes – urban (combined with transportation), agricultural, and desert plots (Figure 1). Cross-correlation matrices with canopy cover from 200 survey plots were constructed for each group. After examining these matrices the most correlated images were selected. We used linearity unmixed vegetation fraction image derived from Landsat 03/18/00 for urban plots, unmixed vegetation fraction image from 05/21/00 for desert plots, and NDVI image from 03/18/00 for agricultural plots.

Preliminary analyses of the data and examination of linear regression assumptions were done in SAS program. Where necessary transformations were applied to the dependent variable (observed cover). Then several alternative regression models were computed using the Fortran program called SLOPES (Isobe, 1990). The program computes standard Ordinary Least Squares (OLS), the inverse OLS (Y|X), Bisector OLS (OLS<sub>bisector</sub>), Reduced Major Axis (RMA), Orthogonal, and Mean OLS regression. We conducted model uncertainty analysis and comparison of the first four models. Regressions with combination of better predictive capability and less error were used to create the map (Figure 3).

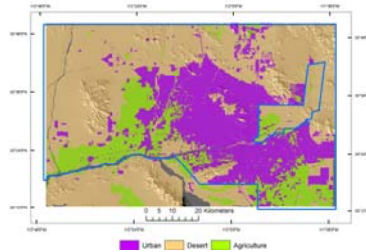


Figure 1. Land use map (reclassified from MAG 2000 land categories)

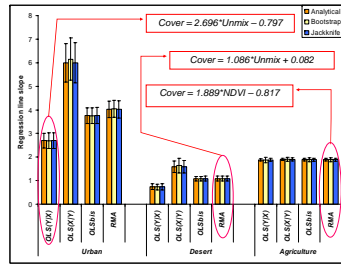


Figure 2. Regression slopes with uncertainties (error bars are  $\pm 1$  standard deviation). Models used in creating the map are shown in red boxes.

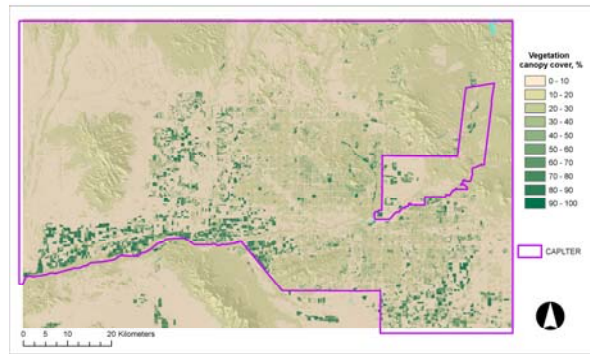


Figure 3. Composite map of vegetation cover in CAP-LTER (March-May 2000). Regression models from Figure 2 are applied selectively for each land use type.

### Acknowledgements

We thank Diane Hope and Corinna Gries for help with 200-point data processing and discussions on its use. Valuable remote sensing data processing suggestions from Dr. Stefanor are appreciated and regression analysis discussions with Dr. Young are acknowledged.

Table 1. Summary of cross-validation

Dataset	Regression method	Variance Ratio	Bias	R	RMSE	SE	N
Urban	OLS (Y X)	0.668	-0.032	0.381	0.157	-0.032	60
	OLS <sub>bisector</sub>	0.957	-0.008	0.378	0.168	-0.008	
	RMA	1.033	0.001	0.376	0.174	-0.001	
Desert	OLS (Y X)	0.863	-0.014	0.294	0.086	-0.014	42
	OLS <sub>bisector</sub>	1.020	0.000	0.301	0.185	-0.014	
	RMA	1.026	0.000	0.301	0.096	-0.014	
Agricultural	OLS (Y X)	0.992	-0.002	0.986	0.045	-0.002	18
	OLS <sub>bisector</sub>	0.998	-0.002	0.986	0.045	-0.002	
	RMA	0.998	-0.002	0.986	0.045	-0.002	

\*Regression methods, variance ratio, bias, RMSE, and SE are explained in the corresponding text. R is correlation coefficient, N is the number of samples.

## ACCURACY ASSESSMENT AND ERROR ANALYSIS

We randomly generated 193 validation sites 90X90 meter in size (=9 Landsat pixels) stratified by major land uses. 18 sites, mostly agricultural, were later discarded. Sampling units bigger than one pixel were chosen to minimize possible geometric errors and ascertain that no one plot is at the edge of the land use it represents. Each site was examined to ensure that it did not overlap with 200 survey plots and it was not a mixed land use one. Some of the removed sites were found to be converted into urban land use, others revealed a significant difference in crop development stage. Each site was then segmented into homogenous patches using high-resolution color aerial photography. Easily identifiable vegetation patches were semi-manually digitized and used to compute the total area (see Figure 4). Percent vegetation cover was then compared with cover estimated by regression models and averaged for 9 Landsat pixels. Same types of error measurements as described above were utilized (Table 2). Plots of predicted versus observed vegetation cover are shown on Figure 5.

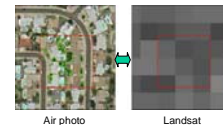


Figure 4. Example of an urban validation site

Table 2. Summary of accuracy assessment

Dataset	Regression method	Variance Ratio	Bias	R <sup>2</sup>	RMSE	SE	N
Urban	OLS (Y X)	1.545	0.077	0.88	0.015	0.077	88
	OLS <sub>bisector</sub>	2.515	0.115	0.87	0.060	0.114	
	RMA	2.790	0.126	0.87	0.080	0.126	
Desert	OLS (Y X)	1.499	0.412	0.39	0.173	0.411	57
	OLS <sub>bisector</sub>	2.341	0.024	0.34	0.010	0.024	
	RMA	2.354	0.024	0.34	0.010	0.024	
Agricultural	OLS (Y X)	1.073	0.012	0.81	0.026	0.012	30
	OLS <sub>bisector</sub>	1.079	0.013	0.81	0.026	0.013	
	RMA	1.079	0.013	0.81	0.026	0.013	

\*Regression methods, variance ratio, bias, RMSE, and SE are explained in the corresponding text. R<sup>2</sup> is coefficient of determination of linear OLS regression constructed with predicted cover being independent variable and observed cover as dependent variable. N is the number of validation sites.

## RESULTS AND CONCLUSIONS

- The effects of different modeling approaches are revealed when models are compared with respect to the statistics summarized in Table 1 and graph in Figure 2. Uncertainties estimated with simulations and cross validation are fairly consistent. Judgments about model performance can also be made based on results of validation using independent samples (Figure 5 and Table 2). Both urban and desert sites are characterized by considerable data scatter, and the difference between slope terms among different models was greater than the errors of any one line (Figure 2, Table 1). All regression methods were biased toward slight underprediction (highest negative bias is associated with OLS) except OLS<sub>bisector</sub> and RMA for desert sites (no bias) and RMA for urban sites (overprediction). OLS<sub>bisector</sub> and RMA by design always exhibit values close to 1 indicating that variance of the observed values is preserved in predicted values. The lowest RMSE is by design associated with OLS however its not very different from RMA constructed for desert plots.
- Overall RMA regression provided a more favorable fit for desert sites (lower slope variability, lower RMSE, closest to unity variance ratio, and no bias), but standard OLS was superior for urban plots (Table 1). However, the use of traditional OLS, which assumes no measurement errors in predictor variables, is considered flawed in principle. We used RMA regression for agricultural plots. Although variability of OLS slope term was somewhat greater than OLS<sub>bisector</sub> and RMA, all three models were essentially identical for agricultural plots.
- Validation results suggest (Table 2, Figure 5) that OLS (Y|X) developed for urban plots is the most accurate model, followed by any one of the three regression of agricultural plots. Desert sites vegetation cover is most accurately predicted by RMA, however this modeling approach it is always possible to have individual predictions outside the range of actual true values (i.e. two fully vegetated urban plots on Figure 5).

## MODEL UNCERTAINTY EVALUATION

Model computations were accompanied by error analysis based on bivariate normal numerical simulations and bootstrap and jackknife resampling of the datasets (available as output of SLOPES program). The two resampling techniques are intended to estimate variability in regression parameters.

Another related model testing procedure, cross-validation, provides a virtually unbiased estimator of prediction error (Efron and Gong 1983). Separate models are developed for each dataset and regression variant by deleting one observation at a time. Each model then is used to predict the observation that was left out. Thus predicted values are compared with actual observed cover. We compared models by computing bias, variance ratios, root mean square errors (RMSE), and standard error (SE) as shown below:

$$Bias = \bar{P} - \bar{P}, \quad Variance\ ratio = \frac{\hat{\sigma}}{\sigma}, \quad RMSE = \sqrt{\frac{\sum_{i=1}^N (\hat{P}_i - P_i)^2}{N}}, \quad SE = \frac{\sum_{i=1}^N (\hat{P}_i - P_i)}{N}$$

where  $\bar{P}$  is the mean of predicted values;  $\bar{P}$  is the mean of observed values;  $\hat{\sigma}$  - standard deviation of predicted values;  $\sigma$  - standard deviation of observed values;  $\hat{P}_i$  - predicted cover for sample  $i$ ;  $P_i$  - observed cover for sample  $i$ ; and  $N$  is the number of observations. RMSE measures the overall accuracy for all samples. Both Bias and SE quantify the effects of systematic errors, such that a positive value signifies overprediction and vice versa. Finally, the variance ratio was used to evaluate how the variance changes with different models. A value close to one would indicate that the variance structure of observed values is preserved in predicted values.

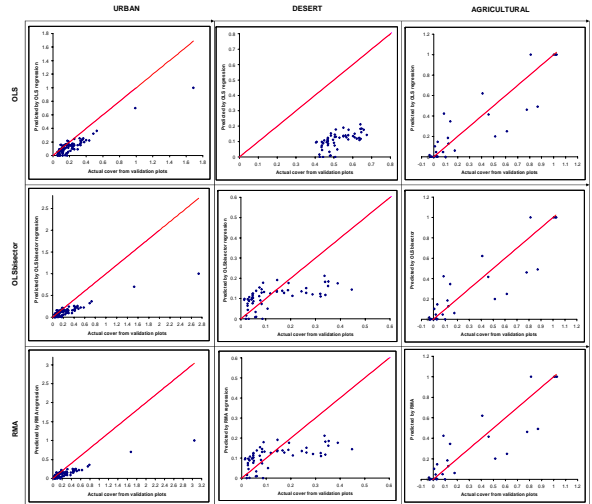


Figure 5. Predicted versus observed (from validation sites) vegetation cover (percent cover). Summary statistics are shown in Table 2.