

**AJAE Appendix:**      **“Eyes in the Sky, Boots on the Ground: Assessing Satellite- and Ground-Based Approaches to Crop Yield Measurement and Analysis”**

**Authors:**              **David B. Lobell, George Azzari, Marshall Burke, Sydney Gourlay, Zhenong Jin, Talip Kilic, and Siobhan Murray**

**Date:**                      **July 25, 2019**

**Note:**                      **The material contained herein is supplementary to the article named in the title and published in the American Journal of Agricultural Economics (AJAE).**

## **Soil Fertility Assessment**

The soil quality index, based on lab analyses of soil samples obtained from the sampled plot locations, is used in our analysis to gauge the possibility of recovering the expected coefficients in production function estimations that use satellite-based yields as dependent variables, as described in greater detail in section 2.3.5. Gourlay et al. (2017) provides details on the collection of soil samples at each plot location in MAPS I. The soil sample collection was not repeated in MAPS II partly due to budget constraints and partly due to the MAPS II preference for the plots that were on the same parcels that had a plot selected in MAPS I, as explained by Gourlay et al. (2017). In MAPS I, four samples of the topsoil (0-20cm) were collected at random locations within each plot and were combined into one composite sample. A single deeper (sub-soil) sample (20-50cm) was collected from the plot center. All samples were shipped to the World Agroforestry Center (ICRAF) Nairobi office, and were subject to spectral soil analysis, with 10 percent of the top- and sub-soil samples also analyzed with conventional wet chemistry testing. The key soil attributes that were measured include pH, texture analysis (sand, % clay, % silt), cation exchange capacity, electrical conductivity (EC), and the concentration of organic carbon (OC), total nitrogen (TN), and potassium.

Following Mukherjee and Lal (2014), a composite soil quality index (SQI) was calculated for each MAPS I plot. Multiple approaches to index construction were employed, including simple additive and weighted additive approaches, as well as a principal component approach and each were computed using topsoil (0-20cm) and subsoil (20-50cm) depths. Bivariate analysis of each index and crop cutting yield estimates (not reported) suggested that the principal component method using top-soil properties was found to correlate more strongly with crop-cut yield than other approaches, and thus, this index is used. Numerous versions of the principal component-based soil quality index were constructed, using different combinations of soil properties. In this approach, principal component analysis (PCA) was first conducted and components with eigenvalues greater than or equal to 1 were retained. Then, the most important variables in each component were identified, including all variables within 10% of the weight of the most important, if the correlation with the most important variable was less than or equal to 0.6. When two or more properties were retained from the same component (where they are weakly correlated and within 10% of the highest weighted property), each property received the same weight.

The index with the greatest predictive power with respect to crop cut yield was composed of organic carbon (%), soil electrical conductivity (an indicator of soil salinity), and pH. These variable values were transformed to a range from 0 to 1, where 1 represents the most optimal value in the sample (e.g., highest value for OC, intermediate values for pH), and 0 represents the lowest value in the sample. A composite index was then generated by weighting each

variable by the fraction of total variance explained by its corresponding component. The relative weights for organic carbon, soil electrical conductivity, and pH are 68.3, 68.3, 31.7, respectively.<sup>1</sup> Given data limitations, the constructed index focuses on nutrient storage capacity but ignores the other two components of soil quality identified by Mukherjee and Lal (2014) related to root development and water storage.<sup>2</sup>

Although these soil samples were acquired in MAPS I, they still provide a useful measure of soil quality to compare with the various yield measures. Importantly, the selected maize plot for most households (n = 312) was part of the same parcel as in the previous year, so that the soil sample was from the same part of their farm. Concerning the remaining sample of households that had a MAPS II plot selected from a non-MAPS I parcel, the median distance between the MAPS II and the MAPS I plot locations was 0.56 kilometers, lending support to likely similarity in soil profiles of nearby plot locations. More importantly, the regression results using soil quality showed very little sensitivity to excluding those households where the parcel moved between years.

---

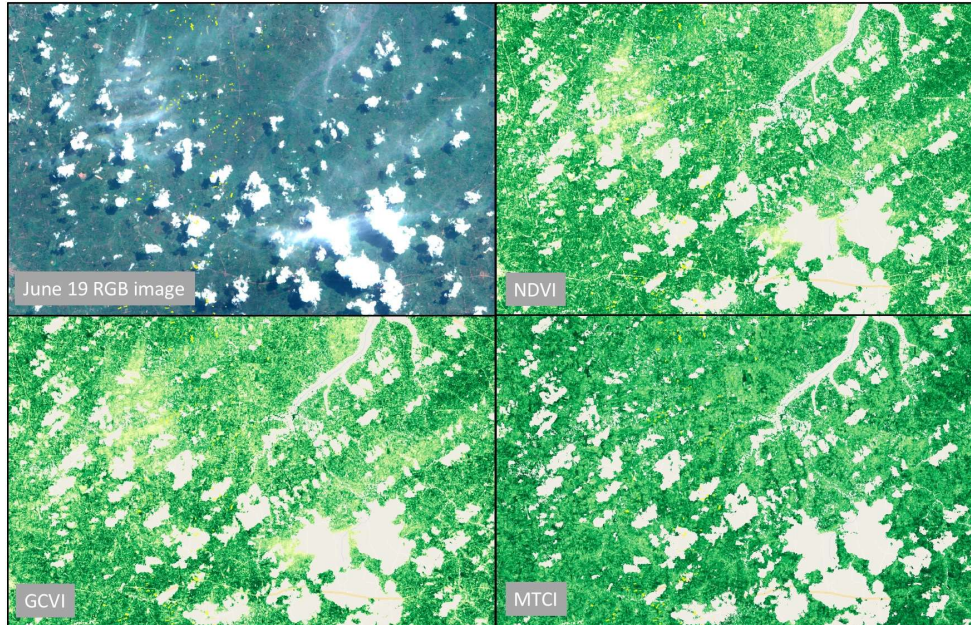
<sup>1</sup> Organic carbon and soil electrical conductivity were both retained from the first component and, therefore, hold the same weight.

<sup>2</sup> The PCA-based soil quality index was constructed for the full MAPS 1 sample, and therefore, analyzes the correlation of soil properties and crop cutting yields on a larger sample than MAPS 2.

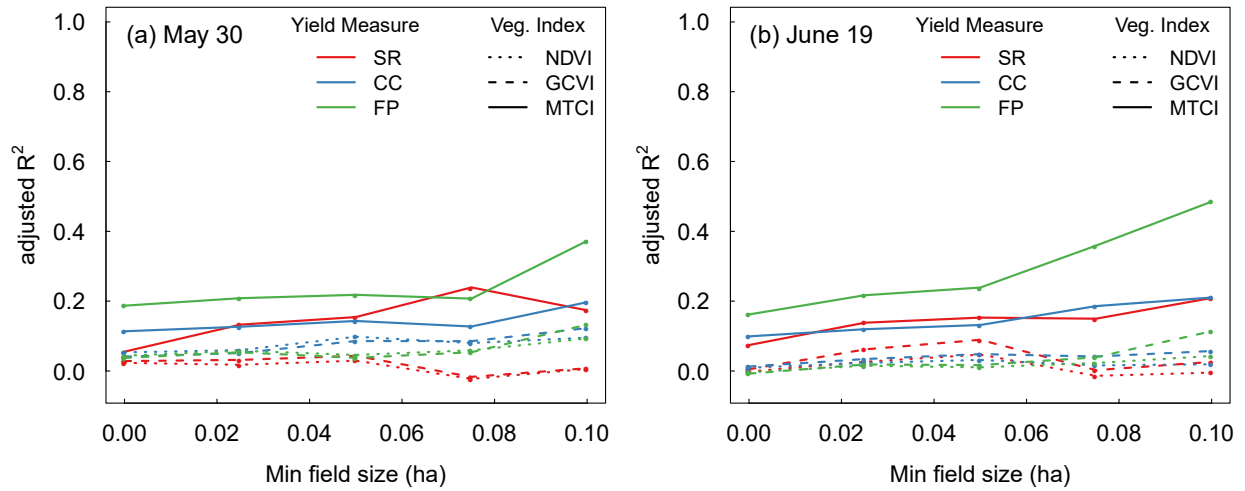
**Table A1. Regression Coefficients for All (Pure Stand + Intercropped) Plots Using Different Yield Measures**

	Dependent Variable/Maize Yield Type					
	Self-report (1)	Crop-cut (2)	Full plot (3)	RS_cal_fp (4)	RS_cal,cc (5)	RS_scym (6)
Log Plot Area (GPS, ha)	-3.37*** (0.47)	0.02 (0.04)	-0.32*** (0.07)	-0.06** (0.03)	-0.0008	-0.0015
Log Plot Distance from Dwelling (GPS, km)	0.21 (0.36)	-0.02 (0.03)	-0.06 (0.05)	-0.06*** (0.02)	-0.04** (0.02)	-0.0006
Cover Crops Present Prior to Planting †	0.18 (0.85)	0.05 (0.07)	0.01 (0.14)	0.01 (0.06)	-0.01 (0.04)	0.01 (0.05)
Log Maize Seed Planting Rate (Kg/Ha)	1.74*** (0.46)	0.03 (0.03)	0.17*** (0.06)	0.04 (0.03)	0.03 (0.02)	0.03 (0.03)
Inorganic Fertilizer Application †	0.70 (1.30)	0.24*** (0.09)	0.34** (0.15)	0.12 (0.07)	0.07 (0.05)	0.07 (0.07)
Log Household Labor Days	0.97** (0.43)	0.01 (0.03)	0.10 (0.06)	-0.03 (0.03)	-0.02 (0.02)	-0.03 (0.02)
Log Hired Labor Days	-0.32 (0.52)	-0.001 (0.03)	0.03 (0.06)	-0.03 (0.03)	-0.02 (0.02)	-0.03 (0.03)
No Hired Labor †	-2.39* (1.22)	-0.09 (0.08)	-0.06 (0.13)	-0.04 (0.07)	-0.01 (0.05)	-0.04 (0.06)
Soil Quality Index	-0.03 (2.84)	0.94*** (0.19)	0.69** (0.34)	0.77*** (0.16)	0.61*** (0.12)	0.68*** (0.14)
Wealth Index	0.13 (0.37)	0.04 (0.03)	-0.06 (0.06)	-0.02 (0.02)	-0.02 (0.02)	-0.02 (0.02)
Agricultural Asset Index	-0.16 (0.37)	0.04 (0.03)	0.07 (0.05)	0.02 (0.02)	0.01 (0.02)	0.002 (0.02)
Dependency Ratio	-0.21 (0.37)	0.02 (0.03)	0.01 (0.04)	0.003 (0.02)	0.004 (0.02)	-0.004 (0.02)
Household Size	-0.10 (0.12)	-0.02* (0.01)	0.005 (0.02)	0.01 (0.01)	0.01 (0.01)	0.01 (0.01)
Manager = Respondent†	0.48 (0.82)	-0.04 (0.07)	0.13 (0.14)	-0.13** (0.06)	-0.09** (0.04)	-0.11** (0.05)
Received Crop-Production Related Extension Services†	-0.01 (0.72)	-0.06 (0.05)	0.02 (0.10)	-0.05 (0.04)	-0.04 (0.03)	-0.03 (0.04)
Female†	0.43 (0.80)	-0.08 (0.06)	-0.04 (0.11)	-0.11** (0.05)	-0.08** (0.04)	-0.0032
Age (Years)	0.01 (0.02)	0.0001 (0.002)	0.01** (0.003)	0.003** (0.001)	0.002** (0.001)	0.003** (0.001)
Years of Education	0.01 (0.08)	-0.002 (0.01)	0.02* (0.01)	0.01 (0.005)	0.005 (0.003)	0.003 (0.004)
Purestand †	-0.21 (0.78)	0.10* (0.06)	0.29*** (0.11)	0.05 (0.05)	0.01 (0.04)	0.01 (0.04)
Log Intercropping Seeding Rate (=100 for Pure stand Plots)	0.07 (0.69)	0.07 (0.05)	0.02 (0.08)	-0.05 (0.04)	-0.0015	-0.05 (0.04)
Constant	-9.65** (4.57)	-0.05 (0.32)	-1.94*** (0.59)	0.42 (0.27)	0.65*** (0.20)	1.89*** (0.24)
Observations	252	463	211	397	397	397
R <sup>2</sup>	0.21	0.14	0.21	0.15	0.15	0.13
Adjusted R <sup>2</sup>	0.14	0.1	0.13	0.11	0.1	0.08
Residual Std. Error	4.96 (df = 231)	0.49 (df = 442)	0.59 (df = 190)	0.38 (df = 376)	0.28 (df = 376)	0.34 (df = 376)
F Statistic	3.07*** (df = 20; 231)	3.55*** (df = 20; 442)	2.57*** (df = 20; 190)	3.33*** (df = 20; 376)	3.25*** (df = 20; 376)	2.73*** (df = 20; 376)

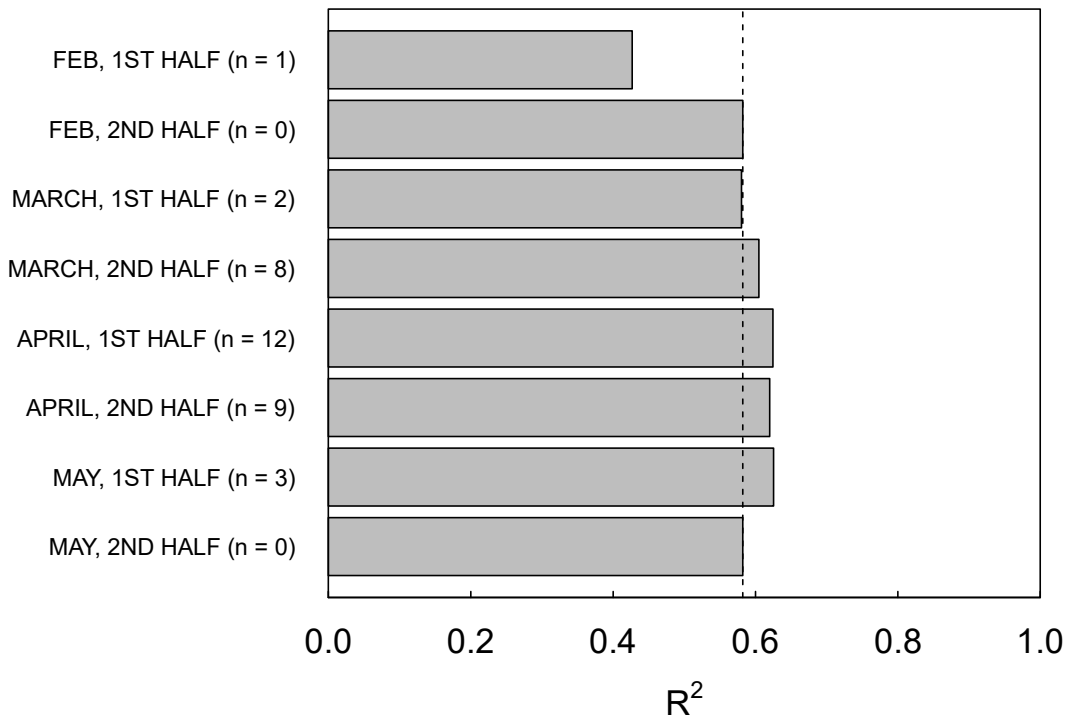
**Notes:** † denotes a dummy variable. \*\*\*/\*\*/\* denote statistical significance at the 1/5/10 percent level, respectively. Standard errors in parentheses.



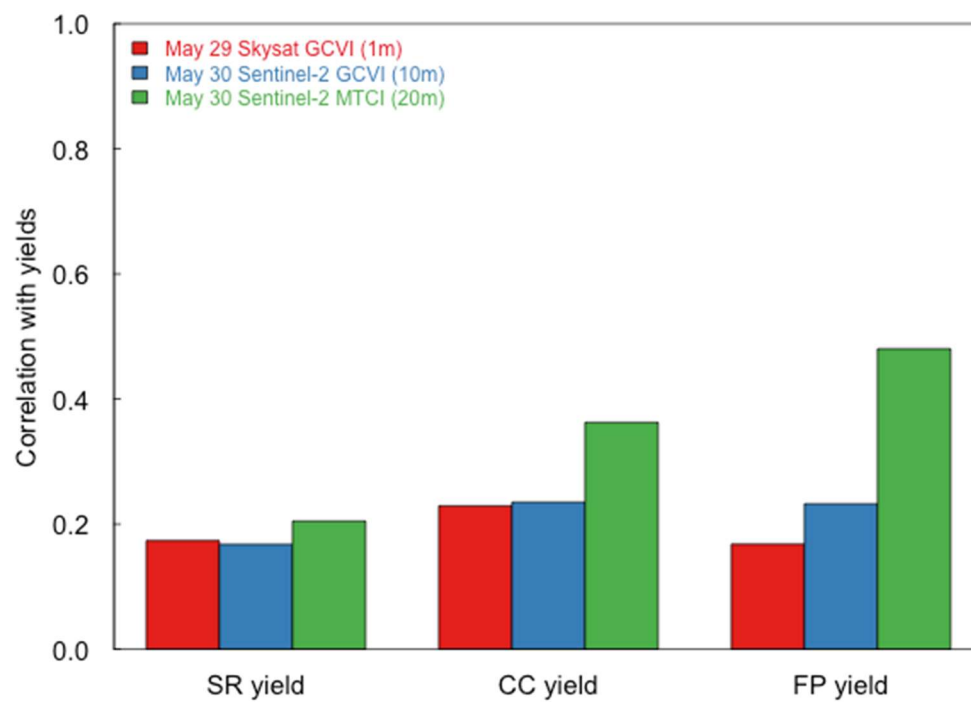
**Figure A1.** The effects of haze on a subsection of the (a) raw red-green-blue reflectance image from June 19, 2016, and the corresponding values of (b) NDVI (c) GCVI and (d) MTCI. For (b)-(d) darker green indicates higher values, and yellow indicates lower values (each VI has a different scale). Areas masked as cloud or cloud shadows are not shown. Both NDVI and GCVI show clear patterns associated with haze, whereas MTCI is less affected.



**Figure A2.** Adjusted  $R^2$  of regressions of yields vs. VI for individual dates, by VI type and type of ground-based yield measure. Models were run for successive subsets of data by excluding plots below indicated plot size. (Same as Figure 3 in main paper but for individual dates). Results for some VIs in Table 2 are not displayed for clarity, but consistently performed worse than GCVI and MTCI.



**Figure A3.** The effect of removing fields with specific sow dates on the agreement ( $R^2$ ) between satellite-based yield estimates and fullplot (FP) yields. The numbers next to each sow date indicate the number of purestand maize fields larger than 0.1 with a FP measurement for that reported sow date. The vertical dashed line shows the  $R^2$  when using all sow dates ( $N = 35$ ).



**Figure A4.** Correlation of different yield measures with VI from Skysat on May 29 or Sentinel-2 on May 30, 2016.