

# Crop Specific Covariate Data based on the NASS Cropland Data Layer for Area Frame Stratification

Claire G. Boryan, Zhengwei Yang  
USDA\NASS\Research and Development Division  
Fairfax, VA

Claire.Boryan@nass.usda.gov  
Zhengwei.Yang@nass.usda.gov



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# Outline

- Background
- Covariate Stratification based on CDL
- Results
- Ultimate Effect Assessment
- Working in Progress
- Conclusions

# Background

- NASS provides timely, accurate, and useful statistics in service to U.S. agriculture

**2001 Wildlife Damage Survey**

**NEWS RELEASE**  
**NATIONAL AGRICULTURE**  
 United States Department of Agriculture  
 Ag Statistics Hotline: (800) 727-6729

**WISCONSIN AGRICULTURAL STATISTICS SERVICE**  
 P.O. Box 8934 Madison, WI 53708-8934  
 In cooperation with WI Department of Agriculture, Trade and Consumer Protection

**2002 Dairy Producer Opinion Survey**  
 November 2002

**Production To Recover**  
 expected to increase in Wisconsin 2 years according to a survey Wisconsin Agriculture Statistics wide survey of producers asked for an assumption that milk prices for the be at the same level as the past five was conducted during May and June

**Based on the survey, 60 percent of producers expect to keep the same herd size, 20 percent plan to increase herd size, and 20 percent intend to discontinue milking by 2007. Actual results will depend on future milk prices, input prices, financing availability, crop yields, and other factors.**

**The number of herds projected for 2007 shows that the diversity of small to large herds will continue. The most prevalent herd size will remain at 50 to 99 cows.**

**Wisconsin Dairy Herds by Herd Size**

Milk cow herd size	May 2002 herds	May 2007 herds (projected) 1/	Change 2007/2002
1-29	2,800	1,440	-45
30-49	4,700	3,440	-27
50-99	7,400	5,600	-24
100-199	1,800	2,080	+9
200-499	700	900	+20
500+	200	440	+120
<b>Total</b>	<b>17,500</b>	<b>13,900</b>	<b>-20</b>

**1/ The May 2007 projection is based on farmers' opinions May-June 2002, with the assumption that milk prices for the next five years will be at the same level as the past five years.**

**My Farmer Plans for May 2007 1/ by Herd Size**

Herds	Keep same herd size	Increase herd size	Discontinue milking
1-29	47	17	36
30-49	71	9	20
50-99	65	19	18
100-199	53	37	10
200-499	33	59	8
500+	22	78	0
<b>Total</b>	<b>60</b>	<b>20</b>	<b>20</b>

**Percent of Herds by Size Group 2007 Projection**

**Legend: Herd Size Groups**

- 1-29
- 30-49
- 50-99
- 100-199
- 200-499
- 500+

**Interactive Data**  
 NASS provides a variety of tools for interacting with our Census datasets.

**Table Lens**  
 Table Lens Application for 1997 Census Data

**Interactive Census Maps**  
 Interactive Census Maps for 2002 Census Highlights

**Interactive Maps**

**Table Lens**

**United States**  
 All data items are from Chapter 2 - Table 1 - Area Summary Highlights: 2002 Selected crops harvested - Land in orchards (acres)

State: United States - County Level | Data Item: Selected crops harvested - Land in orchards (acres)

United States Total: 5,330,439  
 State Total:  
 County Total:

Download data as CSV | XML | PDF

Help | Print | Return to Map

**Legend**

Scale: National

Zero or Data Withheld <= 20,000  
 20,001 to 40,000  
 40,001 to 60,000  
 60,001 to 80,000  
 80,001 to 100,000  
 100,001 >=

Color: Green

Source: USDA-NASS 2002 Census of Agriculture © USDA-NASS 2005-2006

**Navigate:** Mouse-over a specific state/county to view the state/county level data. Right click to zoom (option-click for Mac users). Hold the Alt key and click+drag to pan. For additional assistance with this application, click here to view the support page.

USDA NASS



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# How Does Agricultural Statistics Collected at NASS?

- Agriculture Census every five years
- Estimates from Remote Sensing
- Agricultural surveys
  - Estimates from samples based on NASS area sampling frames(ASFs) and list frames

# What Is An Area Sampling Frames?

- An area sampling frame is a collection of segmented land parcels for the area of interest, such as a state. A land parcel can be defined by its attributes, such as ownership, land usage, land cover, etc.
- NASS ASFs are based on a **stratification of land cover** in the U.S. defined by percent cultivated cropland, i.e. all land parcels are classified different land cover categories!
- NASS Area Sampling Frames have been used as the primary tool to conduct agricultural surveys since 1954.
- The NASS Area Sampling Frames are the basis for the annual June Area Survey in which approximately 11,000 segments are enumerated in early June to collect crop acreage and other agricultural information.

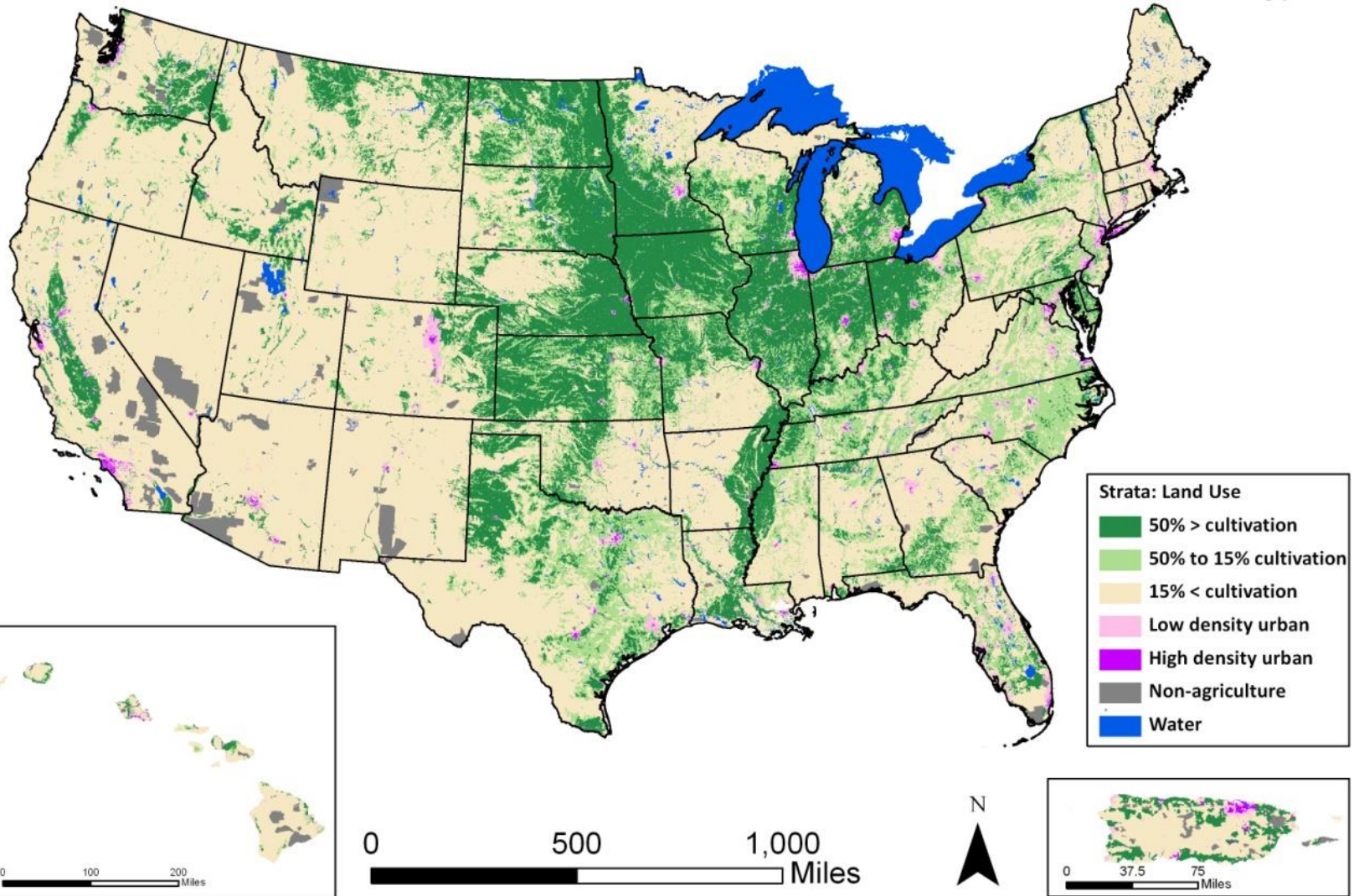


“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”





# National Agricultural Statistics Service Land Use Area Frame



# What Is A Covariate?

- (From Wikipedia) “In statistics, a covariate is a variable that is possibly predictive of the outcome under study. A covariate may be of direct interest or it may be a confounding or interacting variable.”
- For NASS, covariates are variables that may be predictive of where crops will be grown in the future.
- They are derived based on the Area Sampling Frames

# Why Covariates?

- To further improve crop estimates, the ASF Primary Sampling Units (PSUs) have to be substratified based on crop specific information rather than solely on percent cultivation, i.e. the crop specific covariates have to be derived.
- Covariates may be predictive of where crops will be grown in the future.
- The ASFs are built for future use.
- A covariate is a crop specific variable. It can be used to improve area sampling design and ultimately survey estimation for individual major crops.



# How Is Stratification Performed at NASS?

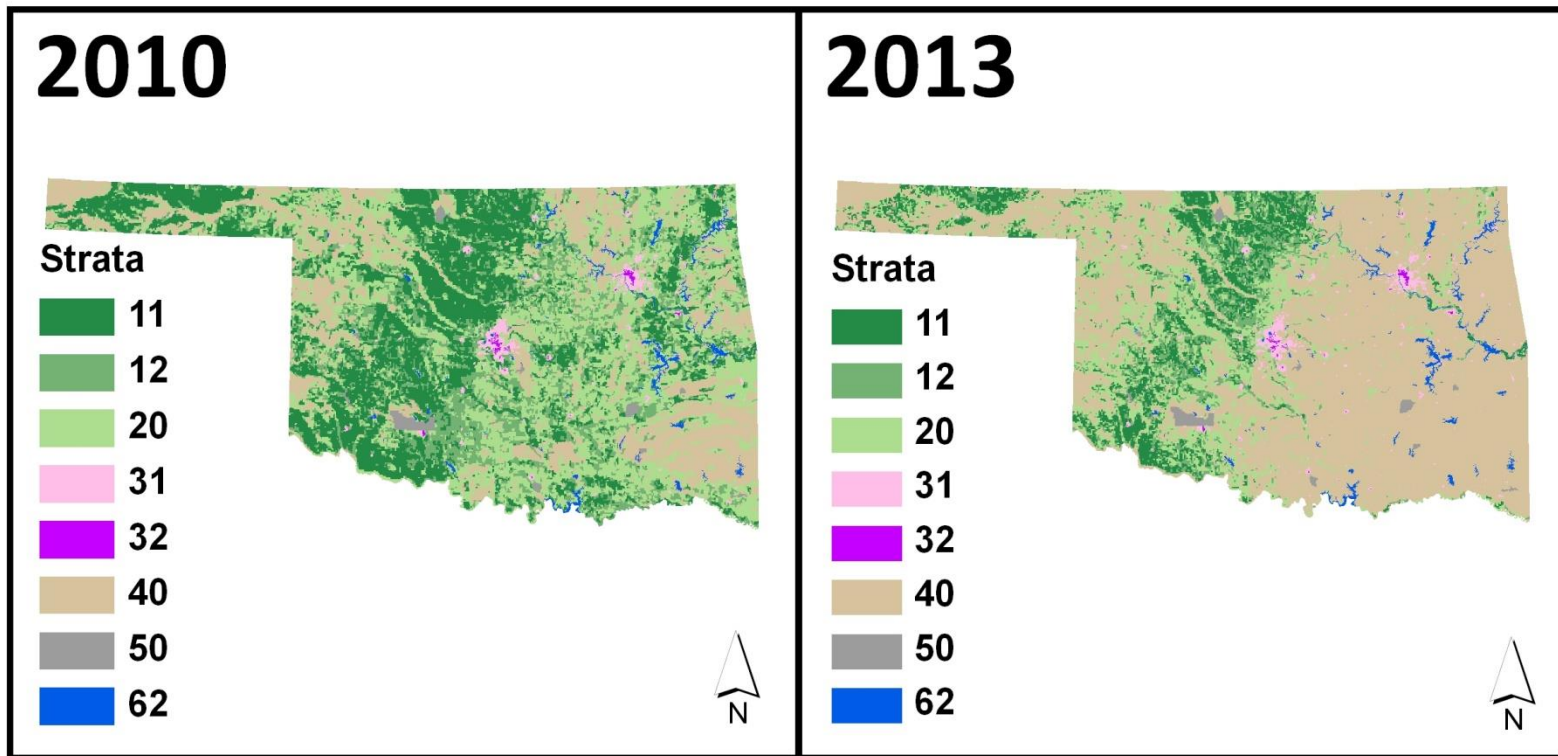
- Stratification goal – make the strata as homogeneous as possible so that the Stratified sampling generally gives more precise (lower variance) estimates for population means and totals than simple random sampling alone.
- Stratification has been conducted by Area Frame staff since 1954 using **visual interpretation** of aerial photography, and later moderate resolution Landsat TM data.
- Digital technology (computerized) was introduced since 1993.
- The NASS Cropland Data Layer products have been used in recent years to aid in the visual interpretation process.
- In the **past two years Cropland Data Layer (CDL) - based automated stratification** has begun to be implemented.



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# CDL Automated Stratification in NASS Operations



**The Oklahoma Area Sampling Frames (2010 and 2013). Stratum 11 (>75% cultivated) was overestimated in the 2010 ASF which was created using the traditional method and updated to more accurately reflect conditions in the 2013 ASF using the CDL automated stratification method. (Graphic courtesy of Kevin Hunt - AF Section -NASS)**

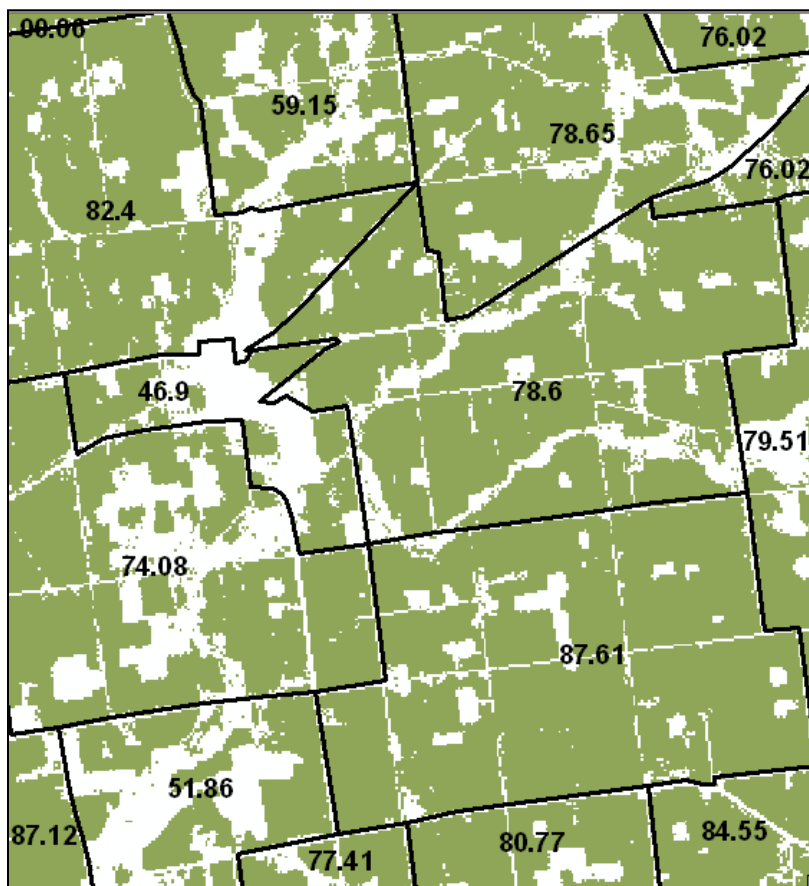
# Automated Covariate Stratification



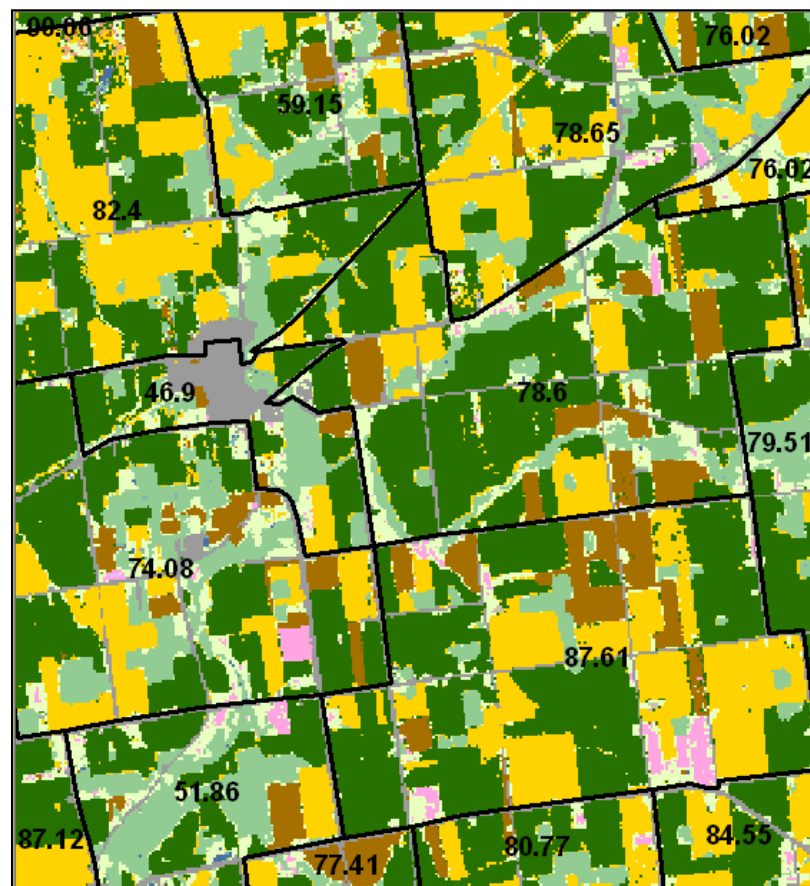
“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# Automated Stratification of the NASS Area Sampling Frame based on the CDL



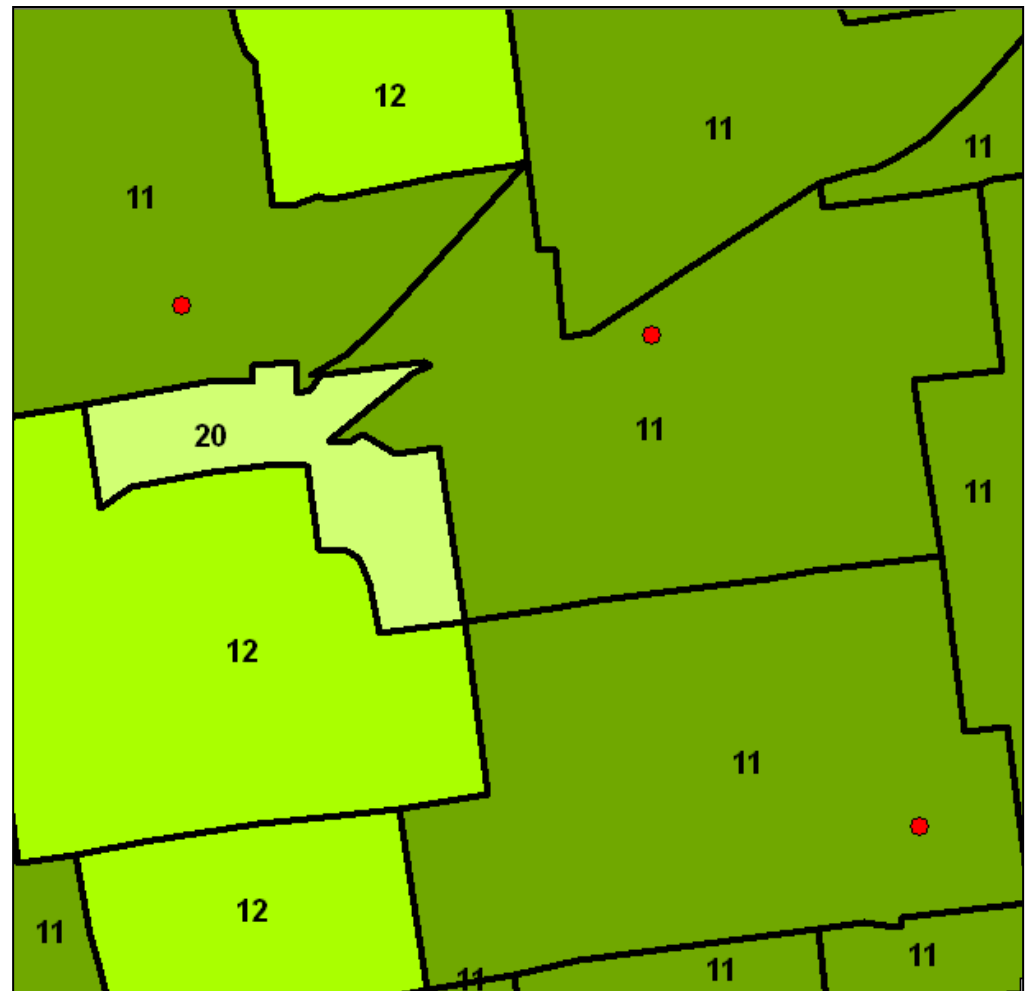
Primary Sampling Units with CDL percent cultivation



Primary Sampling Units with percent crop(s) cover, overlaying a 2010 CDL image product

# CDL based stratification of a NASS Area Sampling Frame (ASF)

- ASF stratification – computing the percent land covered by cultivation (all crops) within a PSU from CDL.
- ASF covariate stratification – computing the percent land covered by a specific crop within a PSU from CDL!
- Classifying each PSU into a defined stratum based on percent covered and the stratum definition.



# Automated Covariate Stratification Procedure

- 1) Derive state level covariate data sets:
  - Select a region of interest from (2007-2010) CDL data, such as a state;
  - Combine the specific crop(s) (i.e. corn/soy, wheat or cotton) from CDL data over multi-year (2007-2010) into one crop category and assigning the corresponding pixels with a value of “1” while grouping the rest of categories into one “other” category and assigning the corresponding pixels with a value of “0”;
  - Save the resulting data into a new covariate data layer ;
- 2) Load and overlay an individual ASF PSU boundary on the CDL covariate data layer;
- 3) Compute percent covariate of each ASF PSU by counting the total number of pixels with value “1” (specific crop) and the total number of all pixels within the PSU boundary. The percent covariate is given by the number of “1” pixels divided by total number of pixels.
- 4) Map each PSU into a defined stratum based on percent covariate covered and the stratum definition.
- 5) For better efficiency, sub-stratification should be conducted:



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# Substratification

- Problem: Find the assignment of  $N_h$  sampling units to  $H$  strata that ***minimizes the sample size***

$$n = \sum_{h=1}^H n_h$$

- Subject To

$$T_j \geq \sum_{h=1}^H \frac{N_h^2 S_{h,j}^2}{n_h}$$

- Where  $T_j$  is a target variance for commodity  $j$  for stratum  $h$  in  $\{1, \dots, H\}$

$$S_{h,j}^2 = (N_h - 1)^{-1} \sum_{i=1}^{N_h} (x_{i,j} - \bar{x}_j)_{h,j}^2$$

# Stratification Results

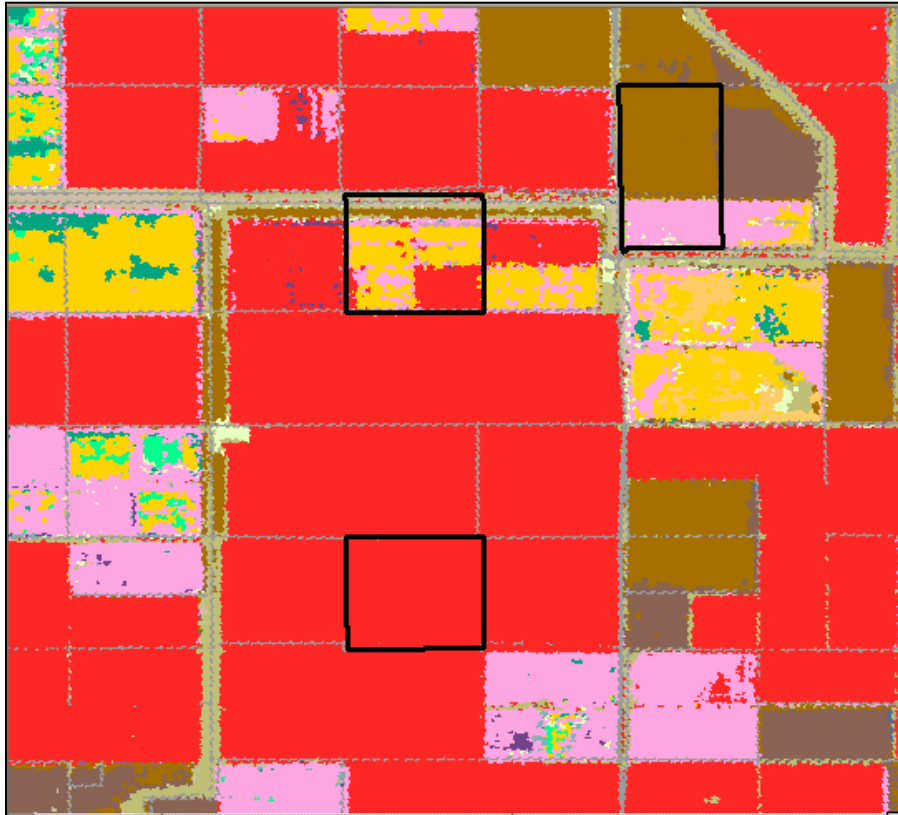


“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”

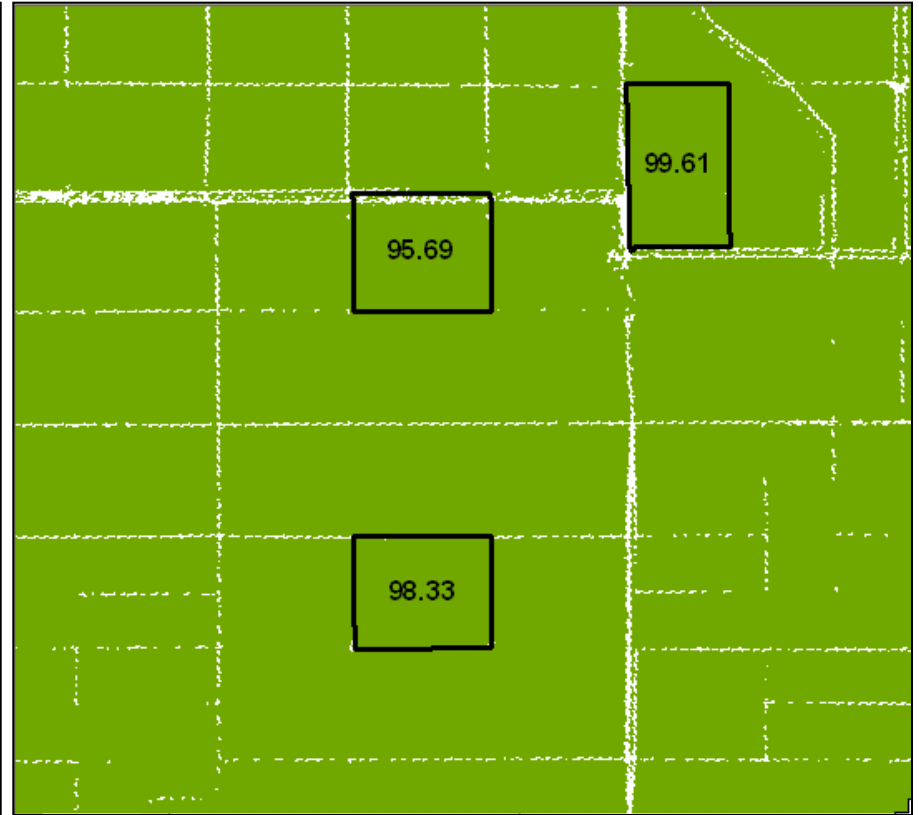




# CDL (2007-2010) covariates at the 2011 June Agricultural Survey (JAS) segment level

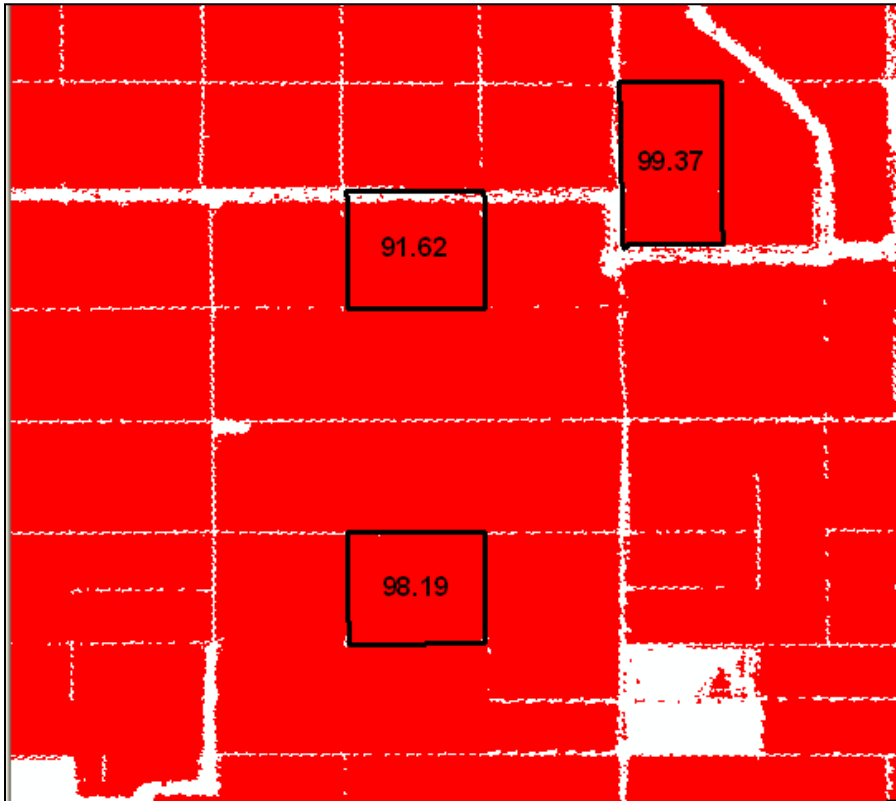


CA11 Cropland Data Layer  
with JAS segments

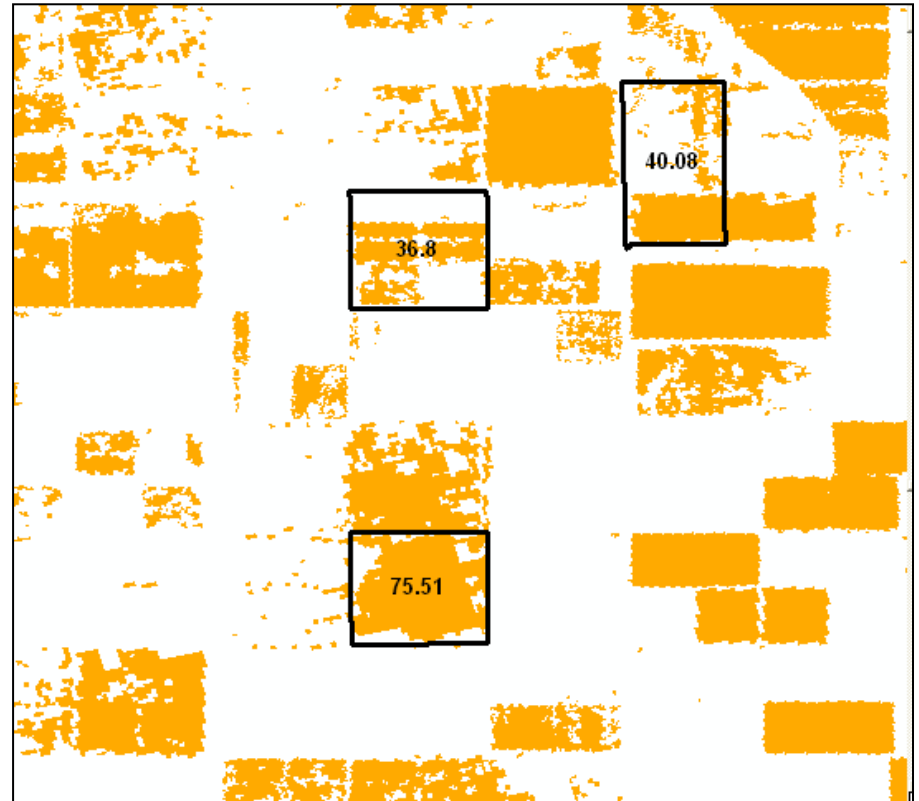


Multi Year (2007-2010) cultivated data set  
with JAS segments (**percent cultivation**  
calculated)

# CDL (2007-2010) covariates at the 2011 June Agricultural Survey (JAS) segment level



Multi-year (2007-2010) **cotton** data set



Multi Year (2007-2010) **corn/soy** data set

# CDL Covariate Predictive Accuracy

	CDL Years	Accuracy	Avg. CDL	Cultivation	Corn/Soy	Wheat	Cotton
California	2007 - 2010	<i>Producer</i>	82.82%	98.95%	52.03%	59.50%	66.73%
		<i>User</i>		95.16%	23.93%	21.06%	36.62%
Indiana	2007 - 2010	<i>Producer</i>	94.82%	96.58%	96.74%	39.88%	N/A
		<i>User</i>		89.08%	86.20%	12.71%	N/A
Mississippi	2007 - 2010	<i>Producer</i>	85.79%	84.11%	93.18%	50.65%	67.55%
		<i>User</i>		93.08%	57.46%	23.08%	36.98%
Nebraska	2007 - 2010	<i>Producer</i>	93.06%	98.45%	94.19%	68.44%	N/A
		<i>User</i>		99.63%	83.76%	25.35%	N/A
Pennsylvania	2008 - 2010	<i>Producer</i>	69.74%	74.16%	83.35%	23.94%	N/A
		<i>User</i>		68.48%	53.11%	8.37%	N/A
Washington	2007 - 2010	<i>Producer</i>	90.27%	89.61%	68.01%	90.04%	N/A
		<i>User</i>		88.78%	27.65%	49.93%	N/A

Validation : 2011 Farm Service Agency Common Land Unit/ NLCD 2006 – cultivation  
 2011 Cropland Data Layers - corn/soy, wheat, cotton

# Ultimate Effect Assessment of Applications of Covariate Data



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# Direct Applications of CDL Covariate Data within NASS

- In the past, commodity information was derived from NASS county level statistics to infer the agricultural makeup for an entire county.
- CDL covariate data sets provide the opportunity to automatically substratify the NASS Area Frame based on commodity information at the Primary Sampling Unit level.



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# Ultimate Assessment - Design Effects

- “The **design effect provides a measure of the precision** gained or lost by use of the more complicated design instead of Simple Random Sampling (SRS)” (*Lohr, S., 2010*)
- The **design effect is defined by the** variance of the estimator from sampling plan (stratified covariate based sampling) divided by the variance of the estimator from a SRS in stratum 11 with the same # of observation units
- **Design effect values less than 1 indicate an increased precision** (reduced variance) in the estimator

# Design Effects

- Comparing prior year design effects using CDL covariate data shows a reasonable overall improvement in substratification efficiency.

Year	Corn	Cotton	Soybeans	Winter Wheat
2012	0.811	0.811	0.773	0.733
2013	0.830	0.683	0.382	0.508

# Conclusion

- The strength of the NASS Cropland Data Layer (CDL) product and CDL based stratification method is **objective and consistent in identification of cultivated cropland.**
- Utilizing the CDL data for Area Frame stratification and sub-stratification will **improve the efficiency, reduce the cost and improve the precision of the June Agricultural Survey estimates.**



# Working in Progress

- 1) Derive covariate data sets from crop planting intensity at each pixel during a multi-year (2008-2012) CDL data, not simple crop coverage.
- 2) Compute percent covariate cover of each ASF PSU with intensity weighting for every pixel.

Thank you

Questions?

[Claire.Boryan@nass.usda.gov](mailto:Claire.Boryan@nass.usda.gov)

[Zhengwei.Yang@nass.usda.gov](mailto:Zhengwei.Yang@nass.usda.gov)

USDA/NASS/RDD

Feb 22, 2012



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



# Procedure for Deriving *Intensity* Covariate Based on Cropland Data Layer (CDL)

- 1) Derive state level covariate data sets from multi-year (2008-2012) CDL data by combining the specific crop (i.e. corn, soy, wheat or cotton) at the pixel level over five year period into five intensity categories and assigning the corresponding pixels with a value of “0,” “1,” “2,” “3,” “4” and “5” indicating the number of years a pixel was planted to the specific crop.
- 2) Load and overlay an individual ASF PSU boundary or grid on the CDL covariate layer;
- 3) Compute percent intensity covariate of each ASF PSU by summing the total number of pixels with values “1 - 5” and the total number of all pixels within the PSU or grid boundary. The percent intensity for a specific crop is given by sum of the pixels divided by total number of pixels.

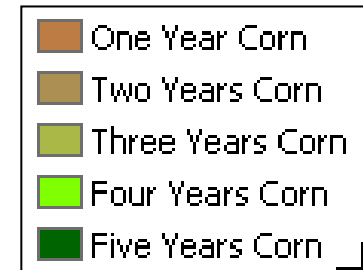
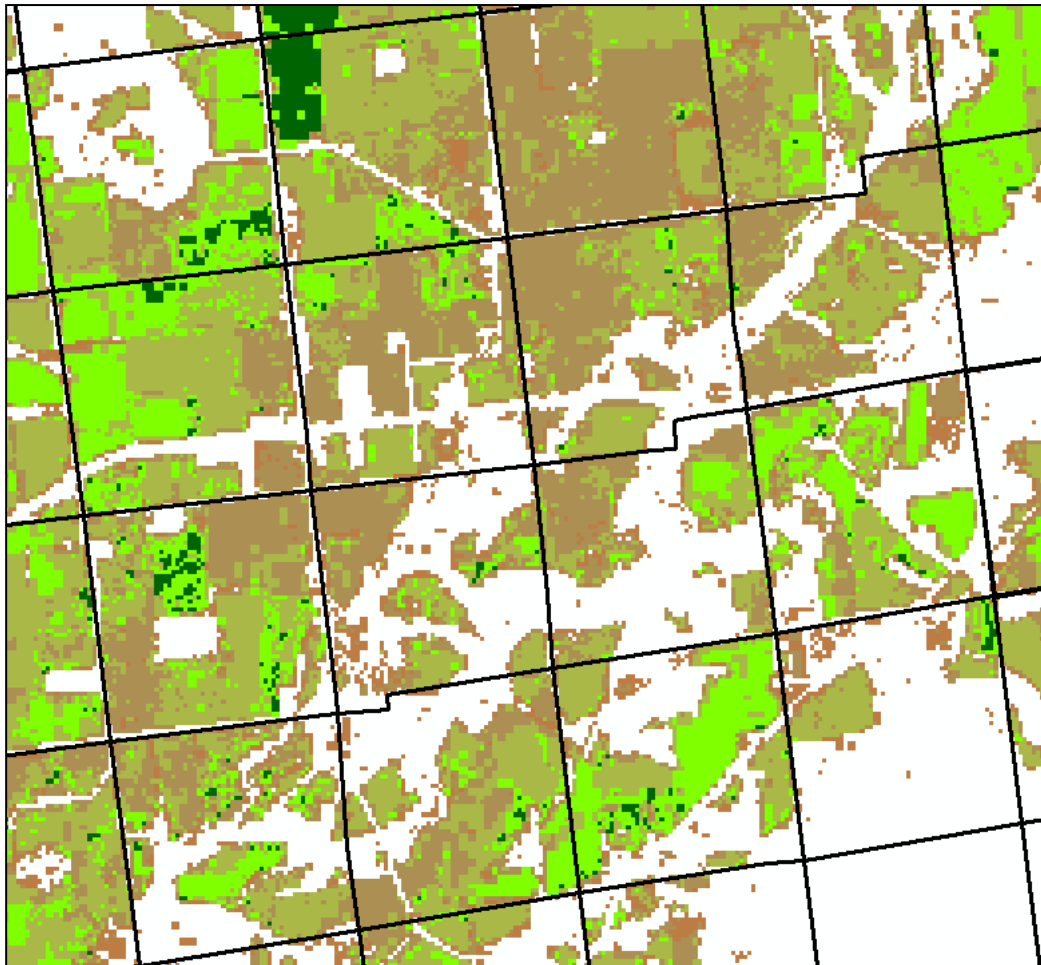
**Comment:** Stratification process is similar to regular covariate stratification. But the stratum definition has to be different.



“... providing timely, accurate, and useful statistics in service to U.S. agriculture.”



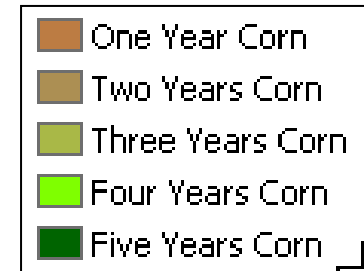
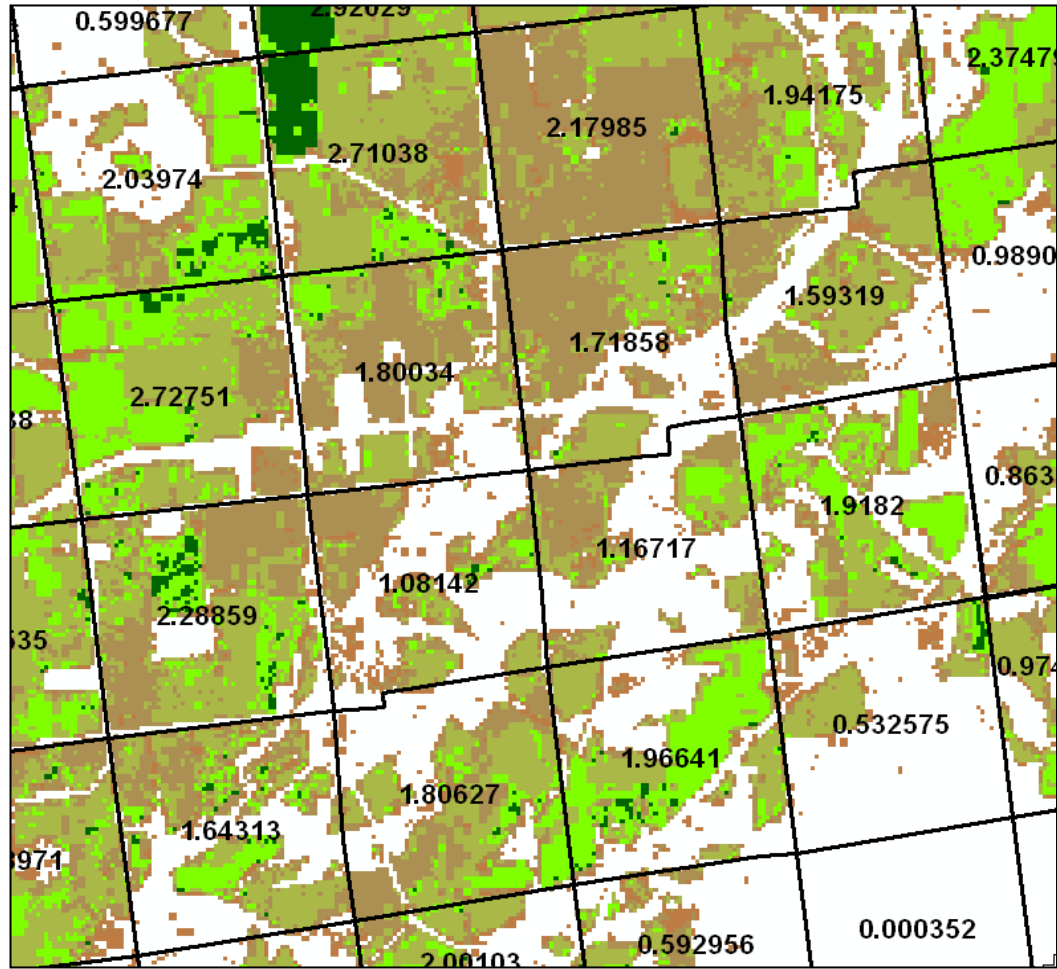
# Indiana CDL (2008-2012) intensity covariates provide improved probability measure



Average corn intensity is calculated at the polygon level.

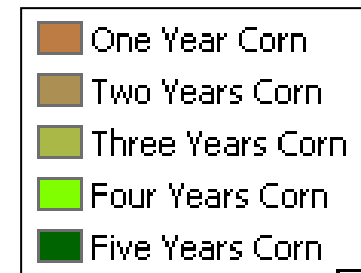
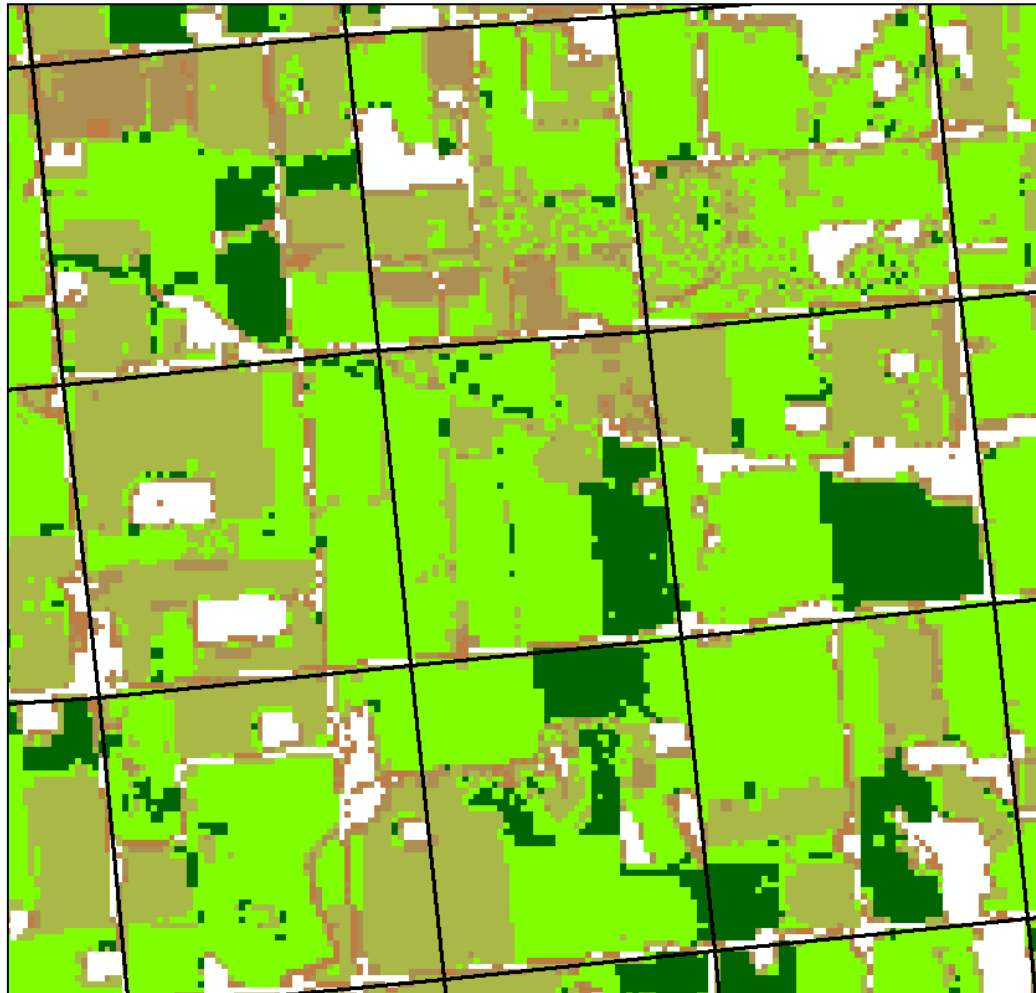
These measurements can be used to more effectively create crop specific clusters for stratification or substratification.

# Indiana CDL (2008-2012) intensity covariates provide improved probability measure



Number of years planted to corn derived at the pixel level (intensity calculation)

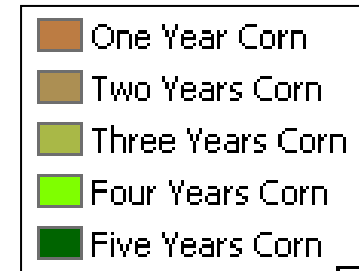
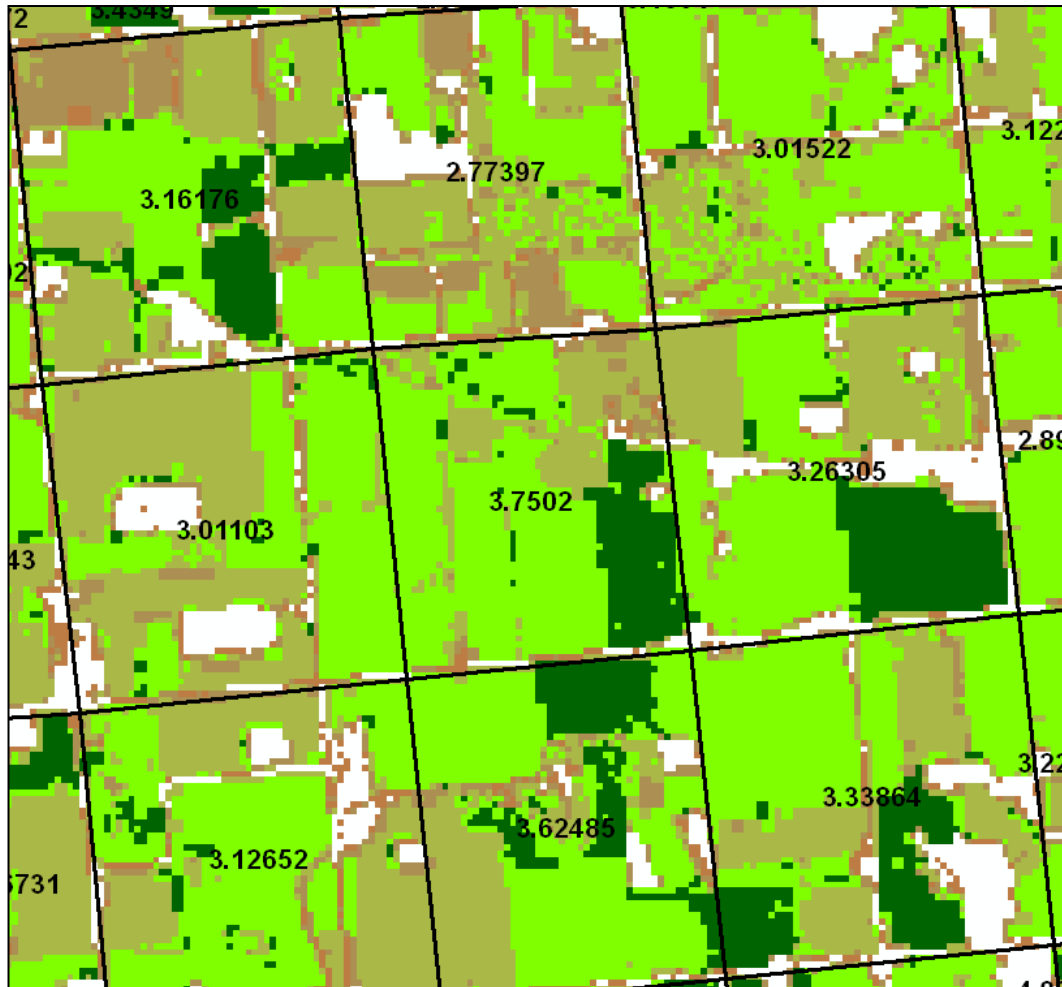
# Indiana CDL (2008-2012) intensity covariates provide improved probability measure



Number of years planted  
to corn derived at the pixel  
level (intensity calculation)

Area of increased corn planting intensity

# Indiana CDL (2008-2012) intensity covariates provide improved probability measure



Number of years planted to corn derived at the pixel level (intensity calculation)

Area of increased corn planting intensity