

Fitting a Bayesian Fay-Herriot Model

Nathan B. Cruze

United States Department of Agriculture
National Agricultural Statistics Service (NASS)
Research and Development Division

Washington, DC
October 25, 2018



Disclaimer

The Findings and Conclusions in This Preliminary Presentation Have Not Been Formally Disseminated by the U.S. Department of Agriculture and Should Not Be Construed to Represent Any Agency Determination or Policy.

Overview

- ▶ NASS interest in small area estimation (SAE)
- ▶ The Fay and Herriot (1979) model
- ▶ Case study: county estimates of planted corn, Illinois 2014
 - ▶ Computation in R and JAGS

Small Area Estimation (SAE) Literature

“A domain is regarded as ‘small’ if the domain-specific sample is not large enough to support [survey] estimates of adequate precision.” –Rao and Molina (2015)

Regression and mixed-modeling approaches in SAE literature

- ▶ Shrinkage–improve estimates with other information
- ▶ Utility of auxiliary data as covariate
- ▶ Variance-bias trade off

Two common models

1. Unit-level models, e.g., Battese et al. (1988)
 - ▶ USDA NASS (formerly SRS) as source of data/funding
2. Area-level models, e.g., Fay and Herriot (1979)



NASS Interest In SAE

Iwig (1996): USDA's involvement in county estimates in 1917

Published estimates used by:

- ▶ Agricultural sector
- ▶ Financial institutions
- ▶ Research institutions
- ▶ Government and USDA

Published estimates used for:

- ▶ County loan rates
- ▶ Crop insurance
- ▶ County-level revenue guarantee

National Academies of Sciences, Engineering, and Medicine (2017)

- ▶ Consensus estimates: Board review of survey and other data
- ▶ Currently published without measures of uncertainty
- ▶ Recommends transition to system of model-based estimates



Fay-Herriot (Area-Level) Model

Fay and Herriot (1979)–improved upon per capita income estimates with following model

$$\hat{\theta}_j = \theta_j + e_j, \quad j = 1, \dots, m \text{ counties} \quad (1)$$

$$\theta_j = \mathbf{x}_j' \boldsymbol{\beta} + u_j \quad (2)$$

Adding Eqs. 1 and 2

$$\hat{\theta}_j = \mathbf{x}_j' \boldsymbol{\beta} + u_j + e_j$$

- ▶ $\hat{\theta}_j$, direct estimate
- ▶ $E(e_j | \theta_j) = 0$
- ▶ $V(e_j | \theta_j) = \hat{\sigma}_j^2$, estimated variance
- ▶ \mathbf{x}_j , known covariates
- ▶ u_j , area random effect
- ▶ $u_j \stackrel{iid}{\sim} (0, \sigma_u^2)$

Fay-Herriot Formulated As Bayesian Hierarchical Model

'Recipe' for hierarchical Bayesian model as in Cressie and Wikle (2011)

Data model:

$$\hat{\theta}_j | \theta_j, \beta \stackrel{ind}{\sim} N(\theta_j, \hat{\sigma}_j^2) \quad (3)$$

Process model:

$$\theta_j | \beta, \sigma_u^2 \stackrel{iid}{\sim} N(\mathbf{x}'_j \beta, \sigma_u^2) \quad (4)$$

Prior distributions on β and σ_u^2

- ▶ Browne and Draper (2006), Gelman (2006): $\sigma_u^2 \sim ?$
- ▶ We will specify $\sigma_u^2 \sim Unif(0, 10^8)$, $\beta \stackrel{iid}{\sim} MVN(\mathbf{0}, 10^6 \mathbf{I})$

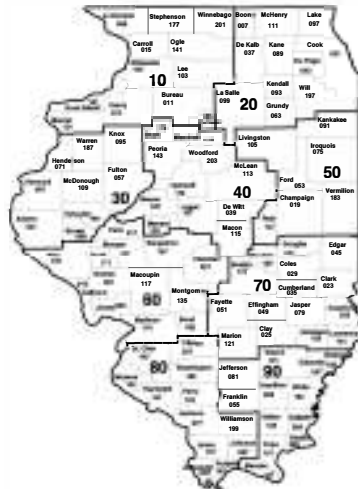
Goal: Obtain posterior summaries about county totals, θ_j

County Agricultural Production Survey (CAPS)

Case study in Cruze et al. (2016)

Illinois planted corn

- ▶ 9 Ag. Statistics Districts
- ▶ 102 counties
- ▶ a major producer of corn
- ▶ End-of-season survey
 - Direct estimates of totals
 - Estimated sampling variances



	Min	Median	Max
n reports	2	47	93
CV (%)	9.1	19.2	92.3

https://www.nass.usda.gov/Charts_and_Maps/Crops_County/indexpdf.php

Covariate x_1 : USDA Farm Service Agency (FSA) Acreage



The screenshot shows the USDA website's 'Crop Acreage Data' page. The page title is 'Crop Acreage Data'. The main content area contains a paragraph explaining that the Farm Service Agency (FSA) provides data on crop acreage for various programs, including the National Conservation Reserve (NCR) and the Conservation Reserve Program (CRP). It also mentions that the data is self-reported by FSA offices and is available for public use. Below the paragraph, there is a section titled 'FSA Crop Acreage Data Reported to FSA' with a list of years: 2011, 2012, 2013, 2014, 2015, and 2016. At the bottom, there is a section titled '2015 Data Year' with a list of links for each year: 2015 (released on 10/13/2016), 2014 (released on 10/13/2015), and 2013 (released on 10/13/2014).

- ▶ FSA administers farm support programs
- ▶ Enrollment popular, not compulsory
- ▶ Data self-reported at FSA office
- ▶ Administrative vs. physical county

<https://www.fsa.usda.gov/news-room/efoia/electronic-reading-room/frequently-requested-information/crop-acreage-data/index>

Covariate x_2 : NOAA Climate Division March Precipitation

Weather as auxiliary variable

- ▶ March: Planting 'intentions'
- ▶ April: Illinois planting
- ▶ **Could rainfall in March affect planting?**
- ▶ One-to-one mapping: ASD and climate division
- ▶ Repeat value for all counties within ASD

	ASD	Precip (in)
	10	1.08
	20	1.35
	30	1.27
	40	1.66
	50	1.50
	60	1.36
	70	1.46
	80	1.69
	90	2.00

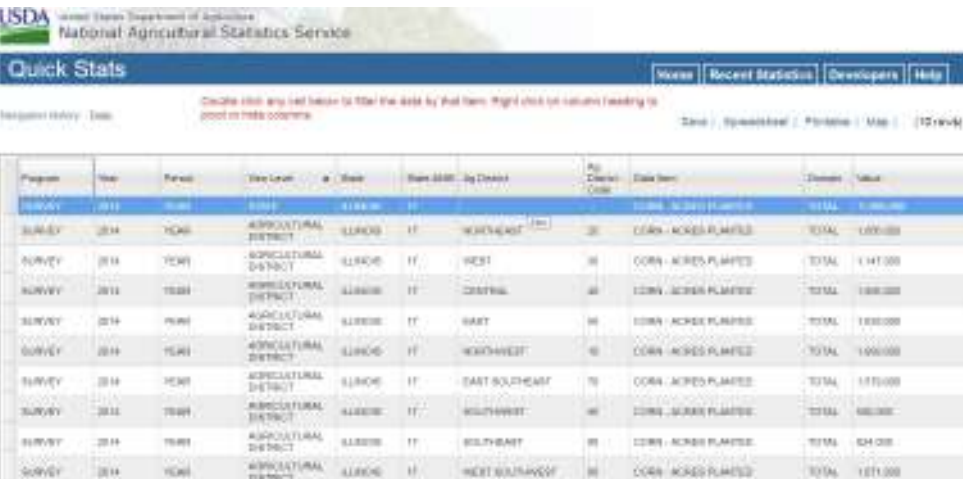
Source: <ftp://ftp.ncdc.noaa.gov/pub/data/cirs/climdiv>
Details in Vose et al. (2014)



NASS Official Statistics

From prior publication: Illinois 2014, 11.9 million acres of corn planted

- ▶ Require: State-ASD-county benchmarking of estimates



The screenshot shows the NASS Quick Stats interface. At the top left is the USDA logo and the text "United States Department of Agriculture National Agricultural Statistics Service". Below this is a "Quick Stats" header with navigation links for "Home", "Recent Statistics", "Developers", and "Help". A search bar contains the text "Corn in any field basis to total the state by dist item. Right click on values leading to pop up table columns." To the right of the search bar are filters for "State", "Year", "District", "Product", and "Map", with "10 rows" indicated. Below the search bar is a table with the following columns: "Program", "Year", "Period", "State", "Dist", "State ASD", "By District", "Ag District Code", "Data Item", "Totals", and "Value". The table contains 10 rows of data for the year 2014, showing corn planting in acres for various districts: NEAR, WEST, CENTRAL, EAST, NORTHWEST, EAST SOUTHEAST, SOUTHWEST, SOUTHEAST, and WEST SOUTHWEST.

Program	Year	Period	State	Dist	State ASD	By District	Ag District Code	Data Item	Totals	Value
SURVEY	2014	YEAR	ILLINOIS	ALL	11,900,000	ALL	00	CORN - ACRES PLANTED	TOTAL	11,900,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	NORTHEAST	20	CORN - ACRES PLANTED	TOTAL	1,000,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	WEST	30	CORN - ACRES PLANTED	TOTAL	1,141,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	CENTRAL	40	CORN - ACRES PLANTED	TOTAL	1,000,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	EAST	50	CORN - ACRES PLANTED	TOTAL	1,000,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	NORTHWEST	60	CORN - ACRES PLANTED	TOTAL	1,000,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	EAST SOUTHEAST	70	CORN - ACRES PLANTED	TOTAL	1,170,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	SOUTHWEST	80	CORN - ACRES PLANTED	TOTAL	960,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	SOUTHEAST	90	CORN - ACRES PLANTED	TOTAL	840,000
SURVEY	2014	YEAR	ILLINOIS	AGRICULTURAL DISTRICT	11,900,000	WEST SOUTHWEST	95	CORN - ACRES PLANTED	TOTAL	1,071,000

State/district: <https://quickstats.nass.usda.gov/results/3A17F375-B762-37BD-8C03-D581DC8F7A85>

County: <https://quickstats.nass.usda.gov/results/478D1A7B-E680-3E5E-95E4-9A59F938A256>

JAGS Model

```
1  ##### Assume this source saved in C:/Your Directory Name/Your_JAGS_model.R
2  model{
3      for(j in 1:m){          #Looping over counties, g=102 for Illinois
4
5          #Defines `data model`--note-JAGS uses precision
6              thetahat[j] ~ dnorm(theta[j], 1/vhat.dir[j])
7
8          #Defines `process model`
9              theta[j] ~ dnorm(beta0+beta1*X1[j]+beta2*X2[j], sigma2u.inv)
10         }
11
12     ## Priors:
13     sigma2u ~ dunif(0, 10^8)
14     sigma2u.inv <- pow(sigma2u, -1)      #Again, precision
15
16     beta0~dnorm(0, .000001)              #Again, precision
17     beta1~dnorm(0, .000001)
18     beta2~dnorm(0, .000001)
19 }
```

- ▶ Note data, process, prior structure from earlier slide
- ▶ Note distributions parameterized in terms of precision
- ▶ Read into R script as stored R source code or as text string

A Pseudo-Code R Script

```
##### Loading some libraries--assumes functioning JAGS installation
1 library(rjags)
2 library(r2jags)
3
4
5 ##### Your data import and wrangling go here
6 ##### we'll actually fit a model scaled by 'size' (a reports)
7 thetahat<-DirInd/Size          ##### Survey Estimate
8 vhat.dbr<-VarDirInd/Size*1     ##### Estimated Survey Variance
9 etc<-rva_norm/size            ##### rva data
10 K2<-test$ppp,3                ##### MGA Match Description
11
12 ##### initialize model
13 set.seed(101); m<-CDE          ##### set seed, define number of countries
14
15 ##### initialize sampler--plausible initial value
16 ##### for sigma0 based on least squares
17 init.sig <- {summary(init.im.coef)$sigma^2}
18
19 ##### Distinguish data inputs and parameters
20 jags.data <- list("thetahat", "vhat.dbr", "K1", "K2", "p")
21 jags.params <- c("theta", "sigma0", "beta0", "beta1", "beta2")
22
23 jags.inits <- function(){list("sigma0" = init.sig)}      ##### function for initial value
24
25 ##### Execute model; assumes JAGS as source code; object returned is an R-list object
26 jags1<-jags(data = jags.data , jags.inits = jags.inits , jags.params = "C:/Your Directory Name/Your_JAGS_model.R",
27             n.chains = 1, a.iter = 10000, n.burnin = 1000)
```



Analysis of JAGS Model Output

Posterior summaries of parameters—based on 3,000 saved iterates

- ▶ Posterior means, standard deviations, quantiles, potential scale reduction factors, effective sample sizes, pD, DIC

```
Inference for Bugs model at C:/Your Directory Name/Your_JAGS_model.R
  3 chains, each with 10000 iterations (first 1000 discarded), n.thin = 9
n.sims = 3000 iterations saved
```

	mu.vect	sd.vect	2.5%	25%	50%	75%	97.5%	Rhat	n.eff
beta0	97.024	205.223	-297.362	-39.365	94.004	235.130	492.579	1.002	1500
beta1	0.865	0.037	0.790	0.841	0.865	0.891	0.937	1.005	830
beta2	-48.553	118.049	-276.194	-126.387	-48.104	28.315	183.179	1.001	2300
sigma2u	20223.038	11544.842	3252.631	11870.939	18247.001	26419.969	47345.031	1.039	84
theta[1]	3399.432	163.965	3083.123	3296.654	3399.326	3505.508	3719.588	1.002	3000
theta[2]	1982.413	153.739	1690.704	1885.191	1977.139	2076.279	2302.119	1.001	3000
theta[3]	2621.446	149.324	2320.691	2525.084	2620.279	2713.351	2925.278	1.001	3000
theta[4]	1296.049	141.511	1014.616	1209.529	1291.823	1383.444	1582.351	1.001	3000
theta[5]	3456.315	157.861	3120.367	3359.261	3458.199	3557.888	3754.838	1.002	1900

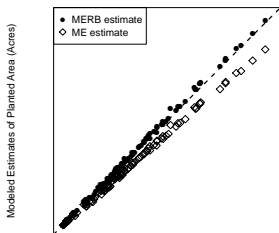
- ▶ Transform back to acreage scale
- ▶ Ratio benchmarking—inject benchmarking factor back into chains as in Erciulescu et al. (2018)



Results: Models With and Without Benchmarking

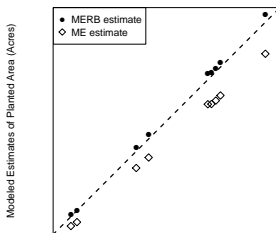
- ▶ Modeled estimates (ME) may not satisfy benchmarking
- ▶ Ratio-benchmarked estimates (MERB) are consistent with state targets and improve agreement with external sources

County Comparisons of Model and FSA Acreage



FSA Planted Area (Acres of Corn)

ASD Comparisons of Model and FSA Acreage

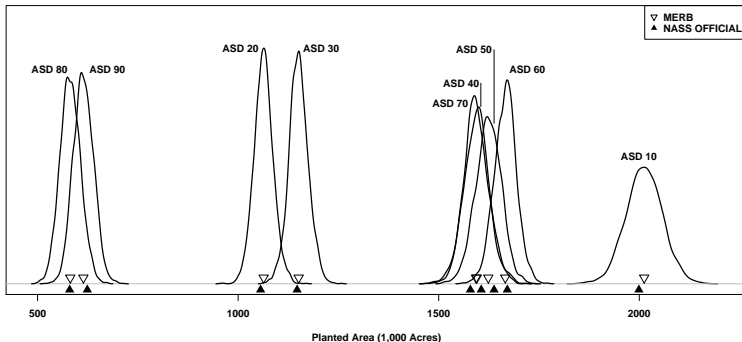


FSA Planted Area (Acres of Corn)

Results: Posterior Distributions of ASD-Level Acreages

Used county-level inputs to produce county-level estimates

- ▶ **Idea:** derive ASD-level estimates from Monte Carlo iterates
- ▶ Sum corresponding draws from county posterior distributions
 - Compute means and variances from aggregated chains



Results: Relative Variability of Survey Versus Model

Obtain estimates and measures of uncertainty for counties and districts

- ▶ Recall the goal of SAE–increased precision!

CV (%) of CAPS Survey Estimates

	<i>Min</i>	<i>Q1</i>	<i>Median</i>	<i>Mean</i>	<i>Q3</i>	<i>Max</i>
County	9.1	16.6	19.2	22.2	23.5	92.3
District	4.4	5.6	6.8	6.6	7.2	8.7

CV (%) of MERB Estimates

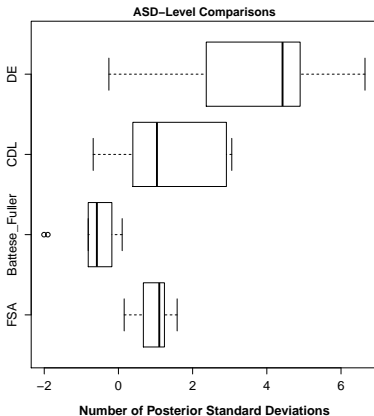
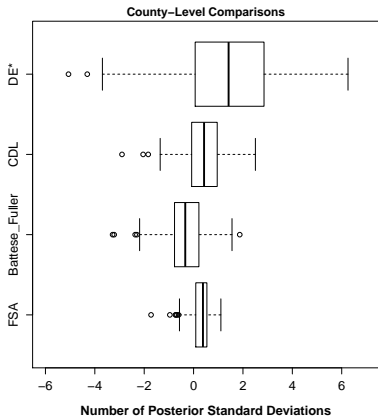
	<i>Min</i>	<i>Q1</i>	<i>Median</i>	<i>Mean</i>	<i>Q3</i>	<i>Max</i>
County	3.6	5.6	7.2	9.0	10.5	31.2
District	1.7	2.0	2.1	2.5	2.3	4.4



Results: Comparison to Other Sources

For counties and districts, compute 'standard score'

- ▶ $(\text{model estimate} - \text{other source}) / \text{model standard error}$
- ▶ Direct Estimates, Cropland Data Layer, Battese-Fuller, FSA



Conclusions

Discussed Bayesian formulation of Fay-Herriot model motivated by NASS applications

Other R packages facilitate Bayesian small area estimation

- ▶ 'BayesSAE' by Chengchun Shi
- ▶ 'hbsae' by Harm Jan Boonstra
- ▶ May be bound by limited choice of prior distributions
- ▶ Transformations of data may be needed

Proc MCMC in SAS added 'Random' statement as of version 9.3

Thanks to Andreea Erciulescu (NISS) and Balgobin Nandram (WPI) for three years of adventures in small area estimation!



References

- Battese, G. E., Harter, R. M., and Fuller, W. A. (1988). An error-components model for prediction of county crop areas using survey and satellite data. *Journal of the American Statistical Association*, 83(401):28–36.
- Browne, W. J. and Draper, D. (2006). A comparison of bayesian and likelihood-based methods for fitting multilevel models. *Bayesian Analysis*, 1(3):473–514.
- Cressie, N. and Wikle, C. (2011). *Statistics for Spatio-Temporal Data*. Wiley, Hoboken, NJ.
- Cruze, N., Erciulescu, A., Nandram, B., Barboza, W., and Young, L. (2016). Developments in Model-Based Estimation of County-Level Agricultural Estimates. In *Proceedings of the Fifth International Congress on Establishment Surveys*. American Statistical Association, Geneva.
- Erciulescu, A. L., Cruze, N. B., and Nandram, B. (2018). Model-based county level crop estimates incorporating auxiliary sources of information. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*. doi:10.1111/rssa.12390.
- Fay, R. E. and Herriot, R. A. (1979). Estimates of income for small places: An application of james-stein procedures to census data. *Journal of the American Statistical Association*, 74(366):269–277.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by browne and draper). *Bayesian Analysis*, 1(3):515–534.
- Iwig, W. (1996). The National Agricultural Statistics Service County Estimates Program. In Schaible, W., editor, *Indirect Estimators in U.S. Federal Programs*, chapter 7, pages 129–144. Springer, New York.
- National Academies of Sciences, Engineering, and Medicine (2017). *Improving Crop Estimates by Integrating Multiple Data Sources*. The National Academies Press, Washington, DC.
- Rao, J. and Molina, I. (2015). *Small Area Estimation*. In Wiley Online Library: Books. Wiley, 2nd edition.
- Vose, R. S., Applequist, S., Squires, M., Durre, I., Menne, M. J., Williams, C. N., Fenimore, C., Gleason, K., and Arndt, D. (2014). Improved Historical Temperature and Precipitation Time Series for U.S. Climate Divisions. *Journal of Applied Meteorology and Climatology*, 53(5):1232–1251.

