

Investigating Covariate Selection for a Bayesian Crop Yield Forecasting Model

Habtamu Benecha

habtamu.benecha@nass.usda.gov

Nathan B. Cruze

nathan.cruze@nass.usda.gov

United States Department of Agriculture
National Agricultural Statistics Service (NASS)

Federal Committee on Statistical Methodology
Research and Policy Conference
March 9, 2018

1/25



Background and Research Questions

- ▶ By mandate, NASS produces monthly crop yield forecasts
- ▶ Official forecasts are consensus estimates of the Agricultural Statistics Board (ASB)
- ▶ Recent research in support of the forecasting program
- ▶ Bayesian hierarchical models
- ▶ Combine data from multiple surveys and covariates

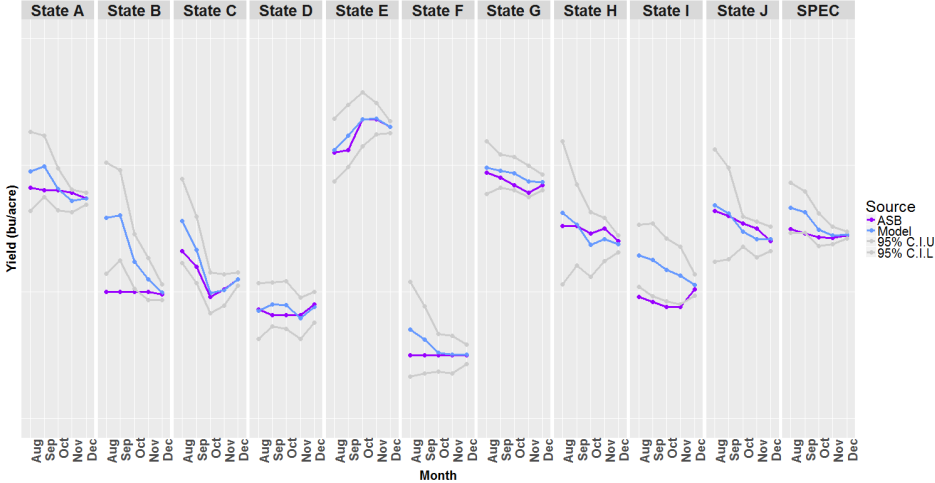
Goal: Which observable covariates are most relevant?

2/25



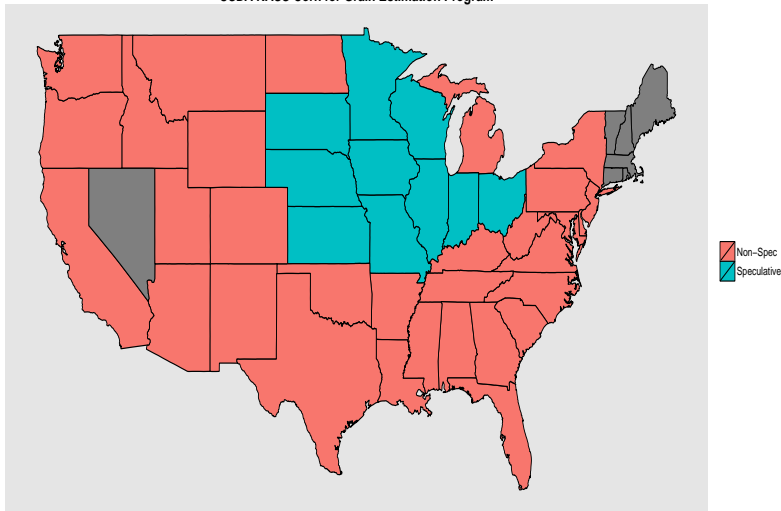
Motivating Example: Forecasting in a Drought Year, 2012 Corn

Corn Yield Forecasts for 2012: ASB and Current Model



Speculative Region for Corn

USDA NASS Corn for Grain Estimation Program



NASS Survey Data and Reporting Timeline

Objective Yield Survey (**OYS**)

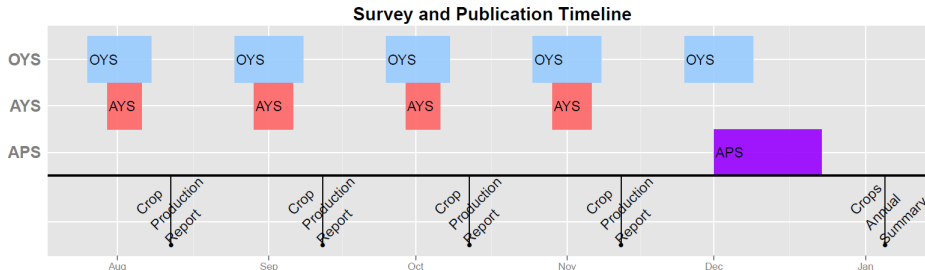
- ▶ Field measurements at sampled plots (Aug.- Dec.)

Agricultural Yield Survey (**AYS**)

- ▶ Interview conducted monthly (Aug.-Nov.)

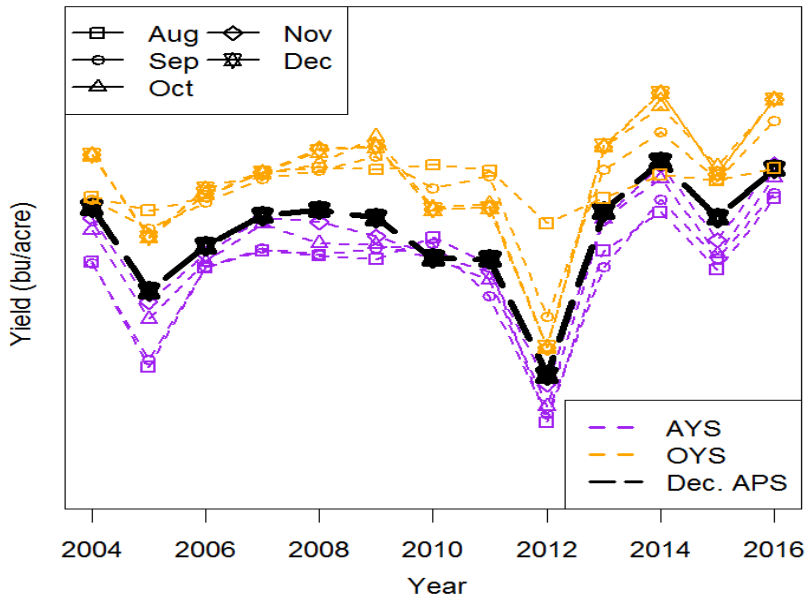
December Crops Acreage, Production, and Stocks Survey (**APS**)

- ▶ Interview conducted post-harvest



Survey Estimates for 2004-2016

Corn Survey Estimates for State A: 2004-2016



Bayesian Hierarchical Model for Speculative Region

Notation

- ▶ μ_t —true yield
- ▶ $t \in \{1, \dots, T\}$ —year index
- ▶ y_{ktm} —observed yield
- ▶ $m \in \{8, 9, 10, 11, 12\}$
- ▶ $k \in \{O, A, Q\}$ —survey index

Stage 1

$$y_{ktm} | \mu_t \sim \text{indep } N(\mu_t + b_{km}, s_{ktm}^2 + \sigma_{km}^2), \quad (1)$$
$$k = O, A; m = 8, 9, 10, 11, 12$$

$$y_{Qt} | \mu_t \sim \text{indep } N(\mu_t, s_{Qt}^2) \quad (2)$$

Stage 2

$$\mu_t \sim \text{indep } N(\mathbf{z}'_t \boldsymbol{\beta}, \sigma_\eta^2) \quad (3)$$

7/25

Bayesian Hierarchical Model for Speculative Region

Diffuse prior distributions on data and process model parameters

- ▶ $\Theta_d \equiv (b_{km}, \sigma_{km}^2)$
- ▶ $\Theta_p \equiv (\beta, \sigma_\eta^2)$

Likelihood function—assuming conditional independence

$$[y_O, y_A, y_Q | \mu_t, \Theta_d] = \prod_{k \in \{O, A, Q\}} [y_k | \mu_t, \Theta_d] \quad (4)$$

Posterior distribution

$$[\mu_t, \Theta_d, \Theta_p | y_O, y_A, y_Q] \propto \prod_{k \in \{O, A, Q\}} [y_k | \mu_t, \Theta_d] [\mu | \Theta_p] [\Theta_d] [\Theta_p] \quad (5)$$

8/25

Bayesian Hierarchical Model–State Level Yield

State-level counterparts indexed by $j \in \{1, 2, \dots, J\}$

Unconstrained State Model–Define $\boldsymbol{\mu}_t \equiv (\mu_{t1}, \mu_{t2}, \dots, \mu_{tJ})$,

$$\boldsymbol{\mu}_t | \mathbf{y}, \Theta_d, \Theta_p \sim \text{indep MVN} \left(\text{vec} \begin{pmatrix} \Delta_{2j} \\ \Delta_{1j} \end{pmatrix}, \text{diag} \left(\frac{1}{\Delta_{1j}} \right) \right) \quad (6)$$

Constrained State Model–Enforce constraint by conditioning state vector $\boldsymbol{\mu}_t$ on $\mu_t = \sum_j w_j \mu_{tj}$

$$(\mu_{t1}, \mu_{t2}, \dots, \mu_{t(J-1)}) \sim \text{MVN}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\Sigma}}) \quad (7)$$

$$\mu_{tJ} = \mu_t - \frac{1}{w_{tJ}} \sum_{j=1}^{J-1} w_{tj} \mu_{tj} \quad (8)$$

9/25

Covariates for the j^{th} State

$$\mu_{tj} \sim N(\mathbf{x}'_{tj}\beta_j, \sigma_\eta^2)$$

Current model for corn includes covariates:

- ▶ T: Trend
- ▶ P: Average July precipitation (NOAA)
- ▶ M: Average July temperature (NOAA)
- ▶ C: Crop condition rating, % rated excellent + good, Week 30 (NASS)

For the Speculative Region: covariate values are defined as weighted averages of state-level covariate values

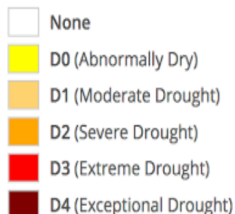
10/25



Additional Covariate

- ▶ Early season model-forecasts
- ▶ Drought severity index

- ▶ $D = \%D3 + \%D4$



- ▶ Pool of available covariates: {T, P, M, C, D}
- ▶ Potential interactions
- ▶ Optimal set of covariates, parsimony

Challenges in Selecting Appropriate Covariates

- ▶ Repeated measure of yield over five months
- ▶ Defining a pool of potential covariates
 - ▶ Crop-specific knowledge
 - ▶ Standard variable selection methods often point to different sets of covariates
 - ▶ Step-wise regression
 - ▶ LASSO
 - ▶ Spike-and-slab regression (Ishwaran and Rao, 2005; Kou and Mallick, 1998), etc
 - ▶ Example: $\{P,M\}$, $\{P,M,C,D\}$, $\{P,M,D\}$ and $\{T,P,M,C,D\}$ - 'best' for the Spec-region in 2016
- ▶ 'Best' sets of covariates depend on state, year and month

Proposed Approach

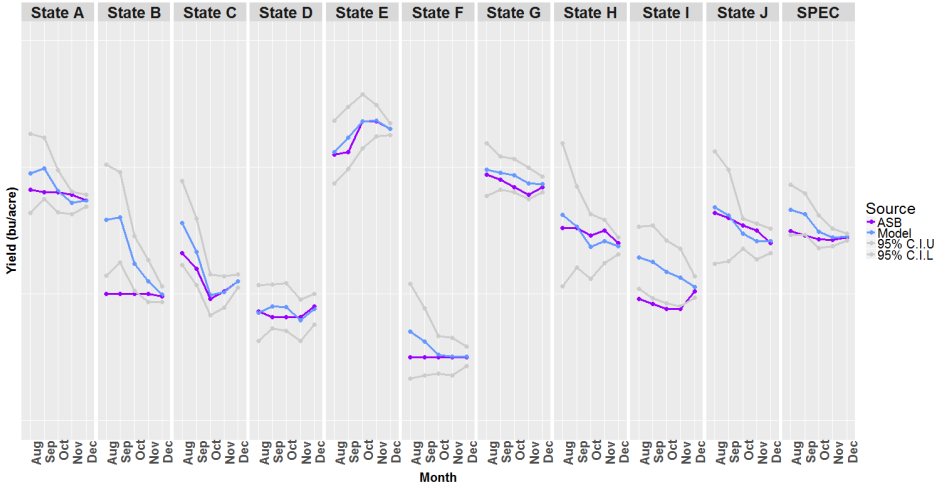
1. Start with alternative sets of covariates that are selected most frequently by traditional variable selection methods
2. Fit models for months from August to December and for years 2012, 2013, 2014, 2015 and 2016.
3. Criteria for decision: percent relative difference from Dec. estimates

$$J = \frac{(\text{Aug. forecast} - \text{Dec. estimate})}{\text{Dec. estimate}} \times 100$$

Model Comparison

- ▶ A total of 17 covariate combinations
- ▶ Subsets of {T,P,M,C,D,TD}
- ▶ Comparisons of models

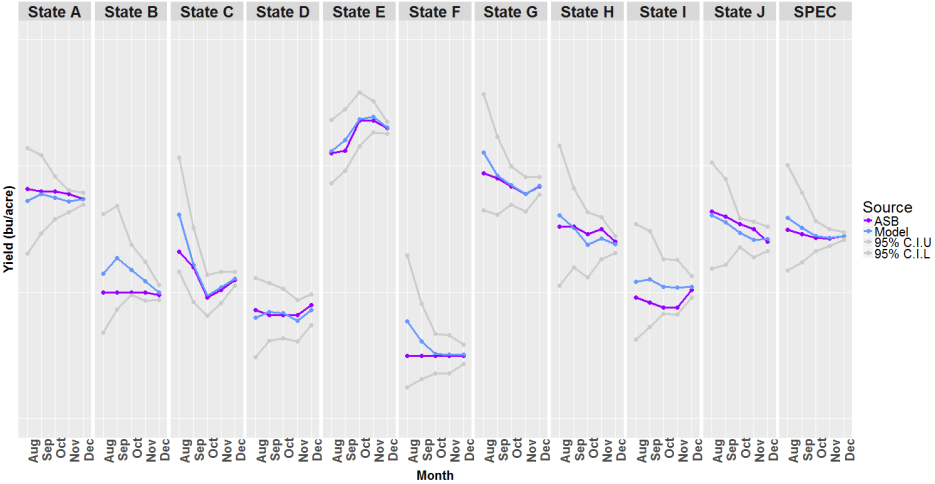
Corn Yield Forecasts for 2012: ASB and Current Model



Model Comparison

- ▶ A total of 17 covariate combinations
- ▶ Subsets of {T,P,M,C,D,TD}
- ▶ Comparisons of models

Corn Yield Forecasts for 2012: ASB and Model {T,P,C,D}

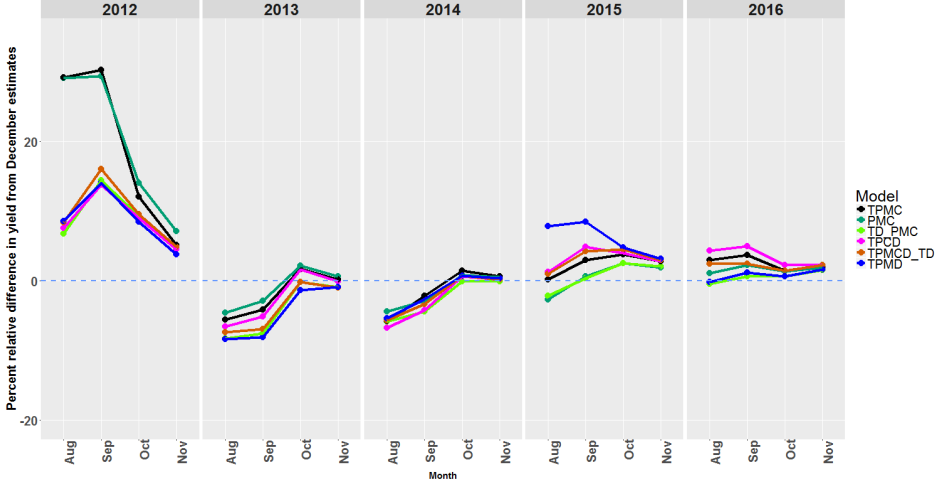


Variables with smallest percent relative differences

State	Covariate-sets		
	Without Drought	With Drought as Main Effect	With Interaction (D*T)
A	-	'13, '15	'12, '14, '16
B	'13, '14, '15	'16	'12
C	'13	'12, '15, '16	'14
D	'12	'13, '14, '15	'16
E	'14	'13, '15, '16	'12
F	'12, '14	'13, '15	'16
G	'12, '13	'14, '15, '16	-
H	'14	'13	'12, '15, '16
I	'12, '15	'13, '14, '16	-
J	'15	'13	'12, '14, '16

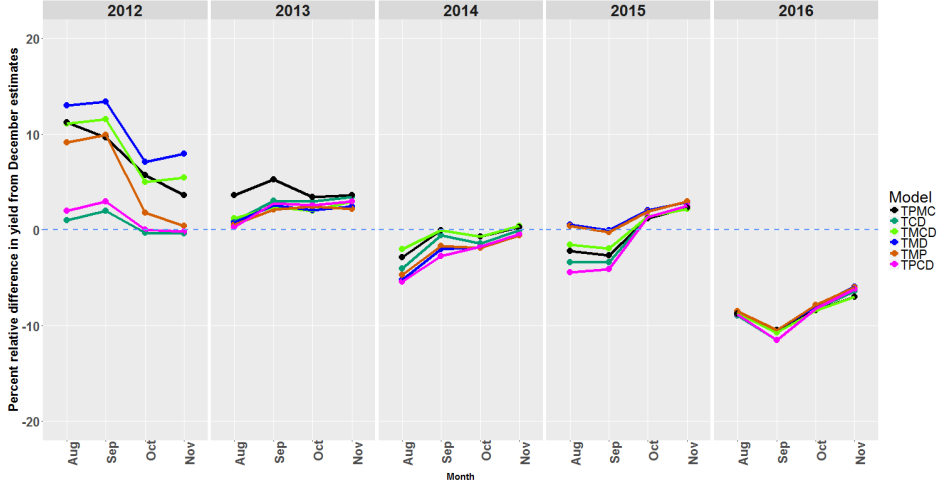
Model Comparison for State B

Percent relative differences from December estimates: State B

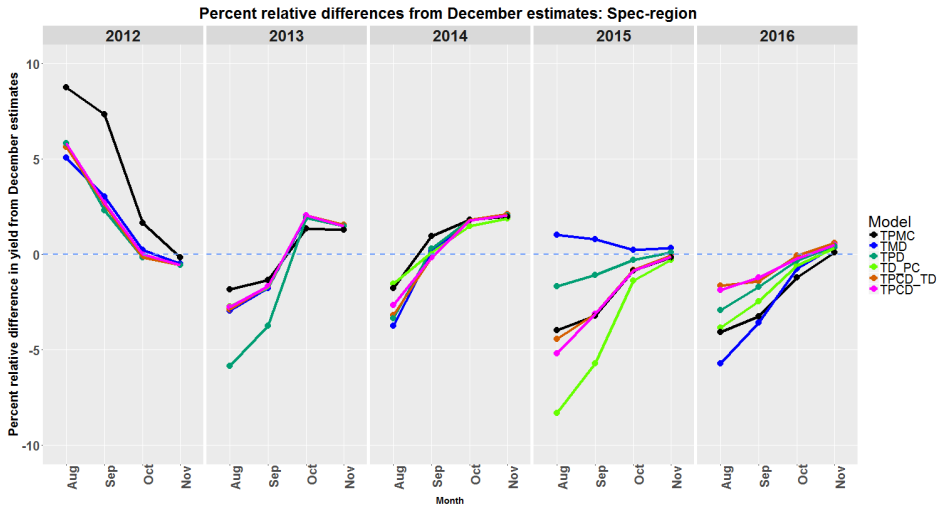


Model Comparison for State I

Percent relative differences from December estimates: State I



Model Comparison for the Spec-region



- ▶ DIC values for December 2016 corresponding to covariate-sets $\{T,P,M,C\}$, $\{T,M,D\}$, $\{T,P,D\}$, $\{P,C, T^*D\}$, $\{T,P,C,D, T^*D\}$, $\{T,P,C,D\}$ are 162.92, 163.06, 163.07, 162.93, 163.15 and 163.02 respectively.

Conclusions

Investigated sensitivity of model forecasts to linear model specification

- ▶ Inclusion of a 'drought' covariate improved early yield forecasts
- ▶ No one-size fits all set of covariates
- ▶ State-specific covariates may be considered

Contact:

`habtamu.benecha@nass.usda.gov`

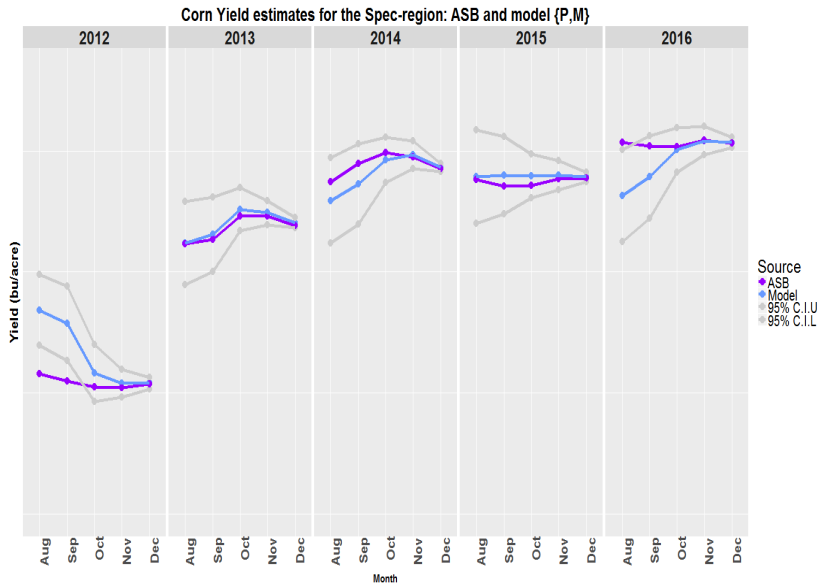
20/25



Select References

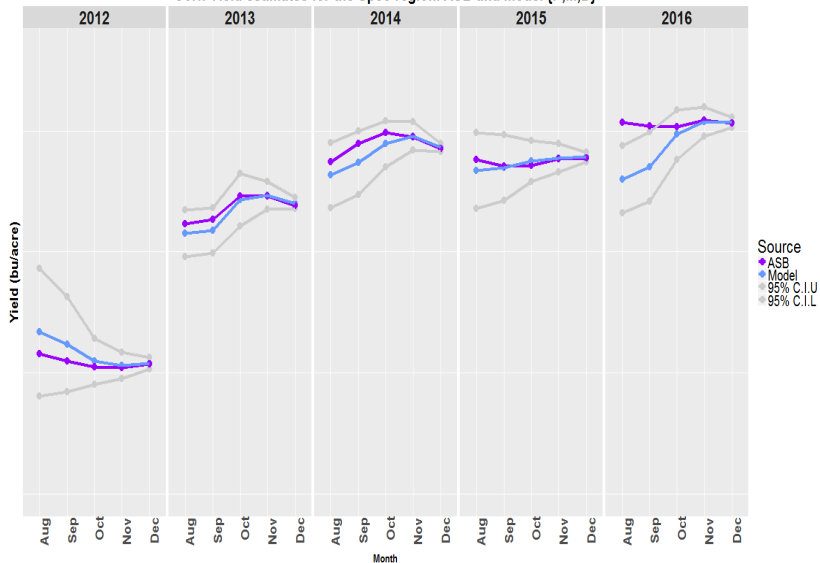
- Adrian, D. (2012). A model-based approach to forecasting corn and soybean yields. Fourth International Conference on Establishment Surveys.
- Cruze, N. B. (2015). Integrating Survey Data with Auxiliary Sources of Information to Estimate Crop Yields. In JSM Proceedings, Survey Research Methods Section. Alexandria, VA: American Statistical Association.
- Cruze, N. B. (2016). A Bayesian Hierarchical Model for Combining Several Crop Yield Indications. In JSM Proceedings, Survey Research Methods Section. Alexandria, VA: American Statistical Association.
- Cruze, N. B. and Benecha, H. (2017). A Model-Based Approach to Crop Yield Forecasting. In JSM Proceedings, Section on Bayesian Statistical Science. Alexandria, VA: American Statistical Association.
- Nandram, B., Berg, E., and Barboza, W. (2014). A hierarchical Bayesian model for forecasting state-level corn yield. *Environmental and Ecological Statistics*, 21(3):507–530.
- Nandram, B. and Sayit, H. (2011). A Bayesian analysis of small area probabilities under a constraint. *Survey Methodology*, 37:137–152.
- Wang, J. C., Holan, S. H., Nandram, B., Barboza, W., Toto, C., and Anderson, E. (2012). A Bayesian approach to estimating agricultural yield based on multiple repeated surveys. *Journal of Agricultural, Biological, and Environmental Statistics*, 17(1):84–106.

Supplementary material 1



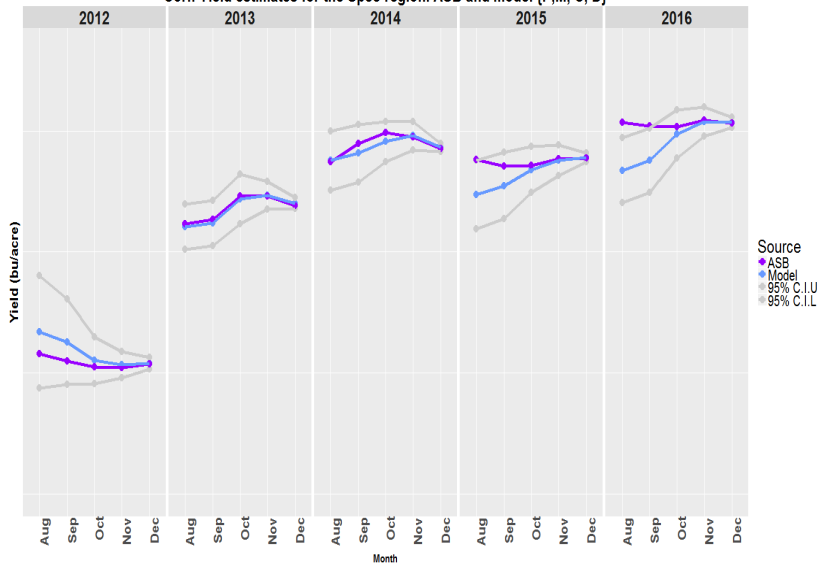
Supplementary material 2

Corn Yield estimates for the Spec-region: ASB and model {P,M,D}



Supplementary material 3

Corn Yield estimates for the Spec-region: ASB and model (P,M, C, D)



Supplementary material 4

Corn Yield estimates for the Spec-region: ASB and model {T,P,M,C,D}

