BAYESIAN STATISTICAL METHODS: ADVANCING FIELD-LEVEL RISK ASSESSMENT IN AGRICULTURE, ACCESSIBLE STATISTICAL TRAINING, AND INCLUSIVE GLOBAL EDUCATION

By

Sarah Manski

A DISSERTATION

Submitted to Michigan State University in partial fulfillment of the requirements for the degree of

Statistics – Doctor of Philosophy

ABSTRACT

Bayesian statistical methods have gained widespread recognition across disciplines due to their intuitive probabilistic nature, incorporation of prior domain knowledge through prior distributions, robust uncertainty quantification, and suitability for working with relatively small datasets. However, the successful implementation, interpretation, and communication of Bayesian methods require a solid understanding of both probability theory and computational techniques. As a Bayesian statistician, I have developed and employed Bayesian methodologies to tackle applied problems across disciplines, collaborating with experts from different fields. Additionally, as a statistics educator, I have designed curriculum to share fundamental skills necessary for comprehending and performing Bayesian analysis.

In this dissertation, I present three projects that illustrate the complexities of utilizing Bayesian methodology in applied problems as a statistician and effectively communicating and teaching the fundamentals of Bayesian theory and application to diverse audiences as a statistics educator. Firstly, I introduce a project that develops Bayesian linear regression and prediction methodologies to quantify the field-level risk mitigation associated with regenerative soil practices in agriculture at a regional scale. Secondly, I discuss the development and execution of an inclusive and accessible workshop aimed at teaching research professionals how to learn the statistical programming language R, as mastering such a language is crucial for practical Bayesian analysis. Finally, I relay how the work from the preceding projects helped build the foundation of a five-day novel training experience to teach the fundamentals of Bayesian statistics to agronomy professionals in Africa.

These projects collectively highlight the multifaceted nature of Bayesian analysis, from its application in addressing real-world challenges to the importance of statistical education and knowledge transfer. By sharing the insights gained through these projects, I aim to contribute to the advancement of Bayesian methodology and facilitate the adoption of Bayesian statistics across disciplines.

CHAPTER 1: INTRODUCTION 1
CHAPTER 2: A BAYESIAN ANALYSIS QUANTIFYING THE RISK MITIGATION EFFECT OF INCREASED ROTATIONAL COMPLEXITY AND WATER STRESS ON CORN YIELD IN THE MIDWESTERN UNITED STATES
CHAPTER 3: "YOU CAN LEARN R": AN ACCESSIBLE AND INCLUSIVE WORKSHOP TO TEACH RESEARCH PROFESSIONALS HOW TO LEARN R
CHAPTER 4: CONCLUSION
BIBLIOGRAPHY

TABLE OF CONTENTS

iv

CHAPTER 1: INTRODUCTION

Bayesian statistical methods are gaining widespread popularity across various disciplines due to their intuitive probabilistic nature, incorporation of prior domain knowledge through prior distributions, robust uncertainty quantification, and ability to handle relatively small datasets within a Bayesian framework. However, employing Bayesian methods effectively requires a solid grasp of probability theory and computational expertise, as researchers need to understand the implementation of these methods, interpret the analysis, and effectively communicate the results.

1.1 Research as a Bayesian Statistician

As a Bayesian statistician, my research interests lie in utilizing Bayesian statistics to address applied problems in interdisciplinary settings through consulting and collaboration. For the past two years, I have been collaborating with an interdisciplinary team of academics and professionals to quantify the risk mitigation effect of regenerative soil health practices, such as increased rotational complexity, on corn yield in the Midwestern United States. Our analysis utilizes observational data spanning 13 years and encompassing over 900,000 fields across 124 counties in Illinois and Minnesota, establishing an empirical connection between increased rotational complexity and reduced risk in terms of corn yield, particularly during periods of water stress. While previous case studies have indicated associations between increased rotational complexity and reduced risk, our analysis is the first of its kind to employ Bayesian methods on such a large scale to quantify this risk mitigation.

By employing Bayesian linear regression, we were able to incorporate domain-specific information about the effects of soil quality, water stress, and rotational complexity on corn yields through prior distributions on model coefficients. To account for spatial heterogeneity and

variation in yield within and between individual fields, we adopted a county neighborhood approach for modeling, incorporating field-level random intercepts and mean-centering for rotational complexity. Our work utilizes the intuitive probabilistic interpretations of posterior predictions provided by Bayesian regression analysis to compare risk probabilities associated with competing management practices across various weather conditions. Consequently, we provide field-level recommendations for adopting more regenerative and stable farming practices.

Two crucial factors in the success of this analysis and effective communication of the methodology and results are the theoretical and computational understanding of Bayesian analysis. Firstly, as a Bayesian statistician, I must possess the ability to interpret and communicate the probabilistic results derived from Bayesian regression analysis. This entails comprehending and comparing posterior predictive probabilities and effectively conveying their interpretation to diverse audiences. Throughout this project, I have summarized and reported Bayesian methodology and results to collaborators, including the founders and strategy officer of the non-profit organization Land Core; agro-ecologist team members at the University of California, Berkeley; federal and MSU grant application reviewers; representatives at the major farm lending cooperative Compeer Financial; and academic audiences through conference talks and manuscripts. In doing so, I tailor explanations of Bayesian methodology and results to each audience based on their statistical knowledge and the understanding necessary for their involvement in the project.

In addition to interpreting and communicating Bayesian results from a theoretical standpoint, conducting Bayesian analysis demands substantial computational understanding and proficiency. In this analysis, every step, from data aggregation and preparation to model fitting,

interpretation, and storage of results, imposes significant computational demands. As the lead statistical analyst on our team, I have acquired the necessary computational skills to interface with our data storage, manipulate data, perform Bayesian analysis, and visualize and store results, all while overseeing peer coding reviews and data checks to ensure the quality and accuracy of our analysis. This project has provided valuable insights into the theoretical knowledge and computational skills required to perform and effectively communicate applied Bayesian analysis.

1.2 Research as a Statistics Educator

The necessity of understanding Bayesian statistics, combined with my passion and experience as a statistics educator, has led me to carefully consider how to teach the fundamentals of Bayesian statistics. As a statistics educator, my goal is to make statistical topics, including Bayesian analysis, inclusive and accessible for researchers in various disciplines. To promote and facilitate the use of statistics among researchers, I have designed and delivered two workshops aimed at making statistical methods more accessible.

An essential aspect of practical Bayesian analysis is the ability to work with a statistical programming language to perform the necessary computational tasks. Therefore, learning a statistical programming language is a fundamental step in utilizing Bayesian methods. Over the course of five years at MSU, I have dedicated a significant portion of my time to learning and teaching techniques in the statistical programming language R. The versatility and extensive data manipulation capabilities provided by R, along with its open-source nature and companion development environment RStudio, make it an ideal tool for researchers at any level. However, I have observed that early-career researchers often learn R in an ad-hoc manner, leading to uncertain learning outcomes.

In the fall of 2022, I was assigned the task of teaching honors introductory statistics to undergraduate students, including guiding them through the steps of statistical analysis using R. Concurrently, I participated in a professional development program focused on facilitating inclusive math learning environments and engaged in a "Teaching as Research" project as part of the Future Academic Scholars in Teaching (FAST) fellowship offered by the graduate school at MSU. Witnessing my honors students' struggles with the steep learning curve of R, combined with my exploration of cooperative learning concepts and universal design in these professional development programs, sparked the idea of designing an accessible and inclusive workshop to teach participants how to learn R. After presenting a cooperative learning exercise on troubleshooting errors in R to the other FAST fellows and observing their interest in learning R for their own research projects, I decided to develop the workshop specifically for early-career researchers like my fellow graduate students. While still in the early stages of developing this workshop, an opportunity arose to use it as a foundational module for a Bayesian workshop targeting agronomy researchers in Africa. I applied for and received the MSU College of Natural Science's Great IDEA (Inclusion, Diversity, Equity, and Accessibility) fellowship to support this endeavor.

In the spring of 2023, I created an R workshop titled "You Can Learn R" with the aim of providing an inclusive and accessible resource for early-career researchers, enhancing their interest, confidence, and ability in learning statistical programming language R. This comprehensive workshop comprises a seminar to introduce participants to the R language, a written document containing curated R learning resources and recommendations, and a supervised working session where participants attempt exercises related to their own research alongside other R learners. A preview and pilot version of this workshop were presented at MSU

to gather feedback and make early adjustments. The first complete iteration of the workshop was conducted in Ethiopia as a foundational module for a larger workshop on the fundamentals of Bayesian statistics in agronomy. This workshop served as a starting point for participants in their R learning journey, offering cooperative exercises and learning materials that they could refer back to as they continue using R for their own research.

In addition to the "You Can Learn R" workshop, I contributed to the development of a larger workshop that provided a unique training experience on the theory and practical application of Bayesian statistics in agronomy to early-career researchers in Africa. This five-day workshop encompassed the "You Can Learn R" workshop, instruction on fundamental probability theory and computational methods required for Bayesian analysis, practical examples of Bayesian analysis in agriculture, and supervised group projects where participants performed preliminary Bayesian analysis on their agronomy data. Drawing on my experience as both a Bayesian statistician and a statistics educator, I facilitated the use of Bayesian statistics for agronomy professionals with minimal R and statistics background.

1.3 Dissertation Outline

In this dissertation, I present three projects that illustrate the complexities of utilizing Bayesian methodology in applied problems as a statistician, as well as effectively communicating and teaching the fundamentals of Bayesian theory and application to diverse audiences as a statistics educator. In Chapter 2, I introduce a project that develops Bayesian linear regression and prediction methodologies to quantify the field-level risk mitigation associated with regenerative soil practices in agriculture on a regional scale. In Chapter 3, I discuss the development and implementation of an inclusive and accessible workshop aimed at teaching research professionals how to learn the statistical programming language R, as

proficiency in such a language is crucial for practical Bayesian analysis. Further, Chapter 3 describes how the work from the preceding projects laid the foundation for a five-day novel training experience aimed at teaching the fundamentals of Bayesian statistics to agronomy professionals in Africa.

These projects collectively highlight the multifaceted nature of Bayesian analysis, from its application in addressing real-world challenges to the significance of statistical education and knowledge transfer. By sharing the insights gained through these projects, I aim to contribute to the advancement of Bayesian methodology and facilitate the broader adoption of Bayesian statistics across disciplines.

CHAPTER 2: A BAYESIAN ANALYSIS QUANTIFYING THE RISK MITIGATION EFFECT OF INCREASED ROTATIONAL COMPLEXITY AND WATER STRESS ON CORN YIELD IN THE MIDWESTERN UNITED STATES

2.1 Introduction

Farming has always been risky, with droughts, floods, heat waves and other hazards harming crop production and farmers' livelihoods for millennia. Climate change increases the severity of these hazards (USGCRP, 2018) including heavier spring rainfall and drier summers (Feng et al., 2016; Swain & Hayhoe, 2015). For example, the 2012 drought reduced maize yields by ~25% in the U.S. Midwest, causing the U.S. government's most expensive year for crop insurance payouts to date, at \$18.6 billion. In 2019, historic spring flooding coupled with summer drought combined to cause a 24% spike in farm bankruptcies over the prior year (Newton, 2019). At the same time, regional specialization in just two crops, maize and soybeans, make the Midwest increasingly vulnerable to stressful weather events (Ortiz-Bobea et al., 2018). While safety nets like crop insurance help mitigate farmers' exposure to negative outcomes, they do not protect food supplies from being disrupted, with concomitant price spikes. Moreover, economic incentives in the current federal crop insurance system encourage simplified crop rotations that may be more vulnerable to stressful weather, and thus can increase risk for the farmer and insurer (Yu et al., 2018).

Several recent syntheses of long-term agricultural experiments show that diversified cropping systems can reduce risks from adverse weather. Previous work using 11 long-term experiments across a continental precipitation gradient in the U.S. and Canada shows that more diverse crop rotations increase corn yields over time and across all growing conditions, including in favorable weather conditions (Bowles et al., 2020). Notably, more diverse rotations also show positive effects on yield under unfavorable weather conditions, with yield losses reduced by 14.0

to 89.9% in drought years. The same pattern holds in seven long-term experiments in Europe where winter and spring cereals had higher yields in diversified rotations as compared with a continuous monoculture (Marini et al., 2020). In particular, yield gains in diverse rotations were up to \sim 1 Mg ha⁻¹ higher in years with high temperature and little precipitation. (Sanford et al., 2021) have further shown that total output from rotations was more stable in rotations with a greater degree of perenniality, and more diverse cropping systems showed less yield decline during drought.

However, all these results are based on plot-level studies from research stations, which do not always translate into field-scale results on working farms (Kravchenko et al., 2017). Thus, how they can be generalized to commercial, working farms remains unclear. Further, climate conditions, intrinsic soil properties, regional management trends, and other factors modulate the extent to which crop rotation promotes resilience to specific stressors in complex, interacting ways, requiring models with widely varying conditions and many data points to sort out yield responses to real-world conditions. Thus, another major knowledge gap is understanding the spatial variation in risk reduction from these practices.

Risk is a function of the severity of a given hazard, the susceptibility or exposure to that hazard, and the response capacity. Risk reduction is a classic example of a benefit that ecosystems provide for people, i.e., an ecosystem service (Wolff et al., 2015), often by reducing susceptibility or increasing response capacity. An example is useful to illustrate this concept, such as the climate risk mitigation associated with the presence of mangroves and tidal marshes. Tidal marshes and mangrove forests can reduce the risk of impacts from coastal flooding on vulnerable communities by moderating flooding impacts (Sheng et al., 2022). But the value of risk reduction from diversifying agroecosystems has rarely been quantified. Crop insurance and

agricultural lending — the two main industries that value risk in agriculture — do not typically reward farmers for changing their production systems to reduce susceptibility to weather and climate hazards. If diversified cropping systems do reduce production and/or profitability risks, then the dollar value of this risk reduction could be applied to insurance and lending policies and passed onto farmers. But achieving these savings will require actuarially sound models that can guide stakeholders in determining the risk reductions associated with these practices in farm-specific contexts.

Only recently have data on agricultural practices and crop yields become available at the field level across wide scales, based on remote sensing and crop modeling (Lobell et al., 2015). Contrary to aggregated data, e.g., at the county level, field-level data allows for understanding fine scale interactions among practices, yields, soils, weather, and other variables. For instance, combining satellite detection of cover crops (Xu et al., 2021) with remotely-sensed and modeled maize yields, Deines et al. (2023) estimated that cover crops reduced maize yields by an average of 5.5% with reduced losses on fields with lower soil ratings, warmer mid-season temperatures, and greater spring rainfall. Importantly, however, prior work leveraging such big data has not examined how diversified cropping systems affect risks from stressful weather, and how this varies over time and space.

In this study, we determined the spatial patterns and magnitudes of cropping system diversification's impacts on rainfed corn yields during stressful dry weather in two contrasting states in the U.S. Corn Belt. We hypothesized that increased rotational complexity would reduce corn yield losses during dry weather without substantial opportunity costs during favorable weather. Using remotely-sensed estimates of crop yields and rotational complexity, we conducted Bayesian statistical modeling focused on corn responses to summer dry periods under

varying levels of rotational complexity. We used data on over 393,000 fields, constructing county-level models to account for spatial variability in yield's response to fixed effects included in the models. We quantified yield responses to adopting crop rotations with more distinct crops and crop turnover, with research questions including the extent to which more complex rotations mitigated the probability of crop yield losses in hot, dry weather, and whether trade-offs exist between benefits in such suboptimal conditions and crop performance under favorable conditions. With climate change expected to increase the frequency and severity of droughts in critical grain producing regions, our results also point toward rotational complexity as an important agricultural climate adaptation strategy. This analysis provides a basis for valuing risk mitigation ecosystem services of complex rotations in actuarial and financial contexts, with implications supporting transitions to more diversified cropping systems.

2.2 Methods

2.2.1 Study Area

Analysis focused on two states in the Midwestern United States, Illinois and Minnesota. Illinois was the second highest corn producer in the country for the entirety of the study period (USDA/NASS, n.d.). Minnesota provides a contrast to Illinois in both geography and rotational practices as a state with distinct climatic constraints and higher average complexity in corn rotations (Socolar et al., 2021).

2.2.2 Data Sources

This study involved processing and aggregating field-level data from a variety of data sources to understand the relationships between corn yield, rotational complexity, soil

characteristics, and weather trends.

Variable	Description	Resolution	Source
Corn yield	Corn yield maps derived from the	30m	(Lobell et al.,
	Scalable Crop Yield Mapper		2015)
Rotational	Value representing the complexity of	30m	(Socolar et al.,
Complexity Index	the crop rotation over the prior six-		2021)
(RCI)	year period, based on the number and		
	turnover of cash crop species		
National Commodity	Proxy for soil quality	30m	gSSURGO,
Crop Productivity			(Natural
Index (NCCPI)			Resources
			Conservation
			Service, 2016)
Minimum and	Indicator of water stress for corn.	4 km	(PRISM
maximum Vapor	Difference (deficit) between the		Climate Group,
pressure deficit	amount of moisture in the air and how		2022)
(VPD)	much moisture the air can hold when		
	it is saturated		
Other weather	Variables such as precipitation,	4 km	Terraclimate,
variables	Palmer Drought Severity Index,		(Abatzoglou et
	Climatic water deficit, temperature,		al., 2018)
	aggregated monthly for the growing		
	season	05 55 1	CL D A C
Soil Moisture	Measures of soil moisture for root	27.75 km	GLDAS,
	zone and for depths including 0 -		(Rodell et al.,
	10cm, 10 - 40cm, 40 - 100cm, and 100		2004)
	- 200cm monthly from May to August		
Soil composition	Available water capacity and average	ırregular	gSSURGO,
	percentage silt, sand, and clay		(Natural
			Resources
			Conservation
			Service, 2016)

Table 2.1. Variables aggregated for inclusion in exploratory data analysis and statistical modeling process, with associated variable description, spatial resolution, and data source.

Yield. Since actual field-level maize yield data are not publicly available, we used maize yield maps derived from the Scalable Crop Yield Mapper (SCYM) (Lobell et al., 2015). The accuracy of SCYM has been extensively evaluated at both the field- and county-scale (Deines et al., 2021; Jin et al., 2017). For instance, when compared with hundreds of thousands of observations from

tractor-based yield monitor data, SCYM showed an R² of 0.45 at the field level, with disagreements likely due to data artifacts in both the yield monitor and satellite sources; when compared with NASS county-level data, SCYM had an R² of 0.69 (Deines et al., 2021). Error in the yield estimates will add noise to our models, but since the SCYM methodology does not include information on crop rotation or other systems-level management, yield estimates from contrasting rotations should not be biased by the algorithm. Compared to other approaches to use remote sensing and modeling to estimate yields, SCYM estimates are the most accurate and widespread (over space and time) dataset on historical corn yields at the field scale in the Corn Belt available (Deines et al., 2021; Kang & Özdoğan, 2019). Other analyses have used SCYM-derived yield maps to evaluate the yield impacts of conservation agriculture practices including two-crop vs. monoculture rotations (Beal Cohen et al., 2019) as well as cover cropping (Seifert et al., 2018) and reduced tillage (Deines et al., 2019).

Rotational Complexity Index (RCI). We calculated RCI for all fields in all years with at least six years of Cropland data layer (CDL) history (the focal year in addition to the five previous) (Boryan et al., 2011). RCI was calculated according to the methods detailed in Socolar et al. (2021). In brief, an index ranging from 0 (least complex) to 5.2 (most complex) was calculated for each field in each year based on the number of crops and frequency of crop turnover in its immediate six-year history. Unlike previous work, RCI was calculated after the cropland data layer history had been aggregated to field level (mode), rather than calculated at the pixel scale and then aggregated to field-level.

National Commodity Crop Productivity Index (NCCPI). We used NCCPI (Albers et al., 2022), a productivity model that ranks the inherent capacity of soils to produce crops without irrigation, as a proxy for soil quality (Li et al., 2016; Seifert et al., 2018; Socolar et al., 2021).

The highest value from the NCCPI submodels (corn/soy, cotton, and small grains) was used for each pixel. NCCPI ranges from 0 to 1 with higher values corresponding to greater soil productivity.

Vapor pressure deficit (VPD). We used VPD as an indicator of corn water stress and agricultural drought, in line with prior studies (Lobell et al., 2014). Vapor pressure deficit is mechanistically linked with corn stomatal regulation and the intensity of water stress (Kimm et al., 2020). Corn yield has a strong negative association with increasing maximum July VPD above a threshold value of ~20 hPa (Xu et al., 2021). We included monthly maximum VPD during the May to August growing season in our exploratory data analysis, and ultimately included maximum July VPD as a model predictor, based on its utility in previous research and the susceptibility of corn to water stress during pollination and flowering.

Gridded corn yields, RCI, and NCCPI were extracted and aggregated to agricultural field level to digitized boundary as constructed in (Yan & Roy, 2016). Field-level time series were constructed by computing the arithmetic mean of values of all grid points that fall within the boundary of each field for each variable. Monthly variables were aggregated to encompass the May to August growing season to examine larger trends in exploratory data analysis.

2.2.3 Exploratory Data Analysis and Spatial Considerations

Due to the breadth of field characteristics and weather variables encompassing 13 years over two states, extensive exploratory data analysis was necessary to understand variability in environmental conditions over space and time. Univariate visualizations such as histograms were analyzed and bivariate correlations and associations between variables were examined, particularly those between our main explanatory variable, corn yield, our primary predictor of interest, RCI, and soil and growing season weather variables known to be primary predictors of

yield. We also examined interactions with rotational complexity and weather variables such as temperature, precipitation, water availability, and soil moisture, as previous case studies have shown that increased rotational complexity is associated with mitigated risk, particularly in stressful weather conditions such as summer drought.

Though general trends between weather conditions and yield were easy to identify, such as a decrease in yield associated with low water availability or high temperatures, variability over space and time made it difficult to identify consistent relationships between rotational complexity and yield or interactions between weather and rotational complexity predicting yield. In a further attempt to model this spatial variation as well as within-field variation over time, we implemented an approach to use the longitudinal nature of our data to separate field-level characteristics from temporal variation. To do this, we partitioned main variables such as RCI and VPD into two parts: a mean value over time for the field and a yearly deviation from that mean. Then to control for spatial variability, we landed on "neighborhood" level models around a focal county. This size was optimal for having a large enough number of fields in a county and its neighbors and encompassing the full variability of a county and its border while being a small enough area to follow similar weather trends.

2.2.4 Bayesian Mixed Effects Model

The Bayesian framework provides a principled way of accounting honestly for model, parameter, and measurement uncertainty quantification. This is critically important when building a predictive model. Bayesian analysis accounts for uncertainty at all levels, with predictions which do not underestimate risk (accuracy), which are not overly pessimistic in their accounting of risk (efficiency), and which typically exceed the predictive power of classical frequentist analysis, particularly in observational studies like ours (see (Dunson, 2001) for

general arguments and (Prost et al., 2008) for the case of yield gaps). In addition, the Bayesian approach has the major benefit of allowing informative priors in the model. In this way our project statisticians and agroecologists could collaborate to elicit meaningful priors based on relationships between explanatory and response variables that are widely accepted in agroecological literature. Furthermore, Bayesian model fitting allows for immediate interpretation of predictions based on posterior distributions as probabilities of future events. This allows us to answer questions about future practices at the field level and compare probabilities of low crop yield under varying weather conditions in order to evaluate risk mitigation. We are then able to quantify these improvements both at the field level and over a larger geographic region to motivate advice on practices for individual farmers as well as institutions with larger interests such as farm lenders.

For each of the included 124 counties in Illinois and Minnesota, we fit a Bayesian mixed effects model to characterize the relationship between field-level yield and crop rotational complexity as defined in the equation in the Figure 2.1, where yields at field *i* in year *t* are modeled as a linear combination of important predictors (elaborated below) with coefficients β_0 , ..., β_8 and an additive field-level random effect α_i , which is zero-centered normally distributed with variance σ_{α}^2 . We also include a normally distributed error term, $\varepsilon_{i,t}$, with variance σ_{ε}^2 .

$$\begin{split} \text{yield}_{i,t} = & \beta_0 + \beta_1 \text{RCI}_{(m)i} + \beta_2 \text{RCI}_{(w)i,t} + \beta_3 \text{VPD}_{i,t} + \beta_4 \text{NCCPI}_i \\ & + \beta_5 \text{VPD}_{i,t} \text{RCI}_{(m)i} + \beta_6 \text{VPD}_{i,t} \text{RCI}_{(w)i,t} \\ & + \beta_7 \text{NCCPI}_i \text{RCI}_{(m)i} + \beta_8 \text{year}_t + \alpha_i + \epsilon_{i,t} \end{split}$$

 $\alpha_i \sim \mathcal{N}(0, \sigma_{\alpha}^2)$ for *i* in 1...*n*

$$\epsilon_{i,t} \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$$

Figure 2.1. Bayesian model formulation for predicting corn yield from RCI, VPD, NCCPI, and year.

In this formulation, we adopt a within-between approach to account for temporal variation within a field as well as spatial variation between fields (Van De Pol & Wright, 2009). We expect that the adoption of higher-complexity crop rotations is confounded with certain unobserved biophysical and socioeconomic factors (e.g., inherent soil quality (Socolar et al., 2021)) in such a way that the causal relationship between RCI and yield may be masked in a simple correlation. We therefore decompose rotational complexity (RCI) into two components: a site mean across all corn years for that field (with subscript (m)) and the deviation between the year's observation and the site mean (with subscript (w) for "within"). For each field, we will call this mean and deviation the "baseline RCI" and "yearly deviation from baseline RCI", respectively. This within-between approach capitalizes on variation in RCI over time within a site. Under this formulation, the effect of the baseline RCI, β_1 , and the baseline RCI interactions with July maximum VPD and NCCPI, β_5 and β_7 , are confounded with any time-invariant unmeasured field attributes, but the effects of the yearly deviation from baseline RCI, β_2 and β_6 , are not confounded with time-invariant values. Coupled with the site-level random intercept, this enables prediction of future conditions that account for inherent field attributes for each field without confounding field-level effects and the effect of RCI.

Based on domain knowledge on the effect of various predictors on corn yield, we define prior distributions for the model coefficients. Based on previous studies, we expect an increase of one unit in RCI to be associated with an increase in yield of 220 - 300 kg ha⁻¹ (Bowles et al., 2020; Seifert et al., 2017). Therefore, for the coefficients for RCI, β_1 and β_2 , we set a normal prior with mean and standard deviation of 260 kg ha⁻¹. Similarly, the detrimental effect of July maximum VPD on corn is approximately linear between 20 and 40 hPa and corresponds to a decrease in yield of 2.2 - 6 Mg ha⁻¹ (Xu et al., 2021). Therefore, an increase of 1 hPa in July

maximum VPD is associated with a decrease in yield of approximately 0.2 Mg ha⁻¹ and we set a normal prior for β_3 with mean and standard deviation of 0.2 Mg ha⁻¹. Soil quality, as represented in our model by NCCPI, is known to have a positive impact on yield, with an increase of 1 in NCCPI corresponding to an increase in yield of 1.1 - 1.2 Mg ha⁻¹ (Deines et al., 2021). We thus set a prior for NCCPI with mean and standard deviation of 1.15 Mg ha⁻¹. Finally, corn yield is known to increase over time due to advances in farming technology, at a rate of approximately 0.15 Mg ha⁻¹ per year (Cassman & Grassini, 2020), so we set a prior for β_8 with a mean and standard deviation equal to that. In each of these cases, we adopt the same mean and standard deviation to create priors that are informative by setting the means based on previous research but also conservative by setting a relatively large standard deviation such that the prior distribution encompasses zero. Due to a lack of domain information on the interactions between variables, the remainder of the model coefficients are fit with default, weakly informative priors as defined by the R package Stan interface, brms, used to fit the models (Bürkner, 2017).

In our Bayesian framework, we fit models by county to accommodate the fact that the relationship between crop practices and yield may vary at large spatial scales. For each county, we fit a model to all data within that county and all adjacent counties (that county's "neighborhood") in order to capture trends in the focal county and its boundary; we then generated predictions and assessed fit for only data from the focal county. The neighborhood modeling approach allows each county's neighbors to support inference on the relationships within that county, minimizing the possibility of edge effects at county borders. Figure 2.2 shows all counties in both Illinois and Minnesota and identifies which counties were included in our modeling and predictions (see Section 2.2.6 on Model Validation and Limitation for details). In Figure 2.2, we highlight a single focal county, Logan County, in blue and its neighborhood in

purple.



Figure 2.2. Map by county showing Illinois and Minnesota counties included and excluded in analysis and a sample focal county with its associated neighborhood.

2.2.5 Interpreting Posterior Predictions to Quantify Risk Mitigation

In each county and for each field, we generate posterior predictive distributions of yield in three weather conditions (July maximum VPD of 18, 20, and 22 hPa representing "normal", "somewhat dry", and "dry" conditions, respectively) for seven levels of rotational complexity. Because the RCI scores also incorporate rates of turnover between crops, there is a range of values that correspond to a given number of crops in rotation. Here we simplify those ranges with six RCI values that are representative of 1-6 crop rotations (Table 2.2), as well as a seventh level that corresponds to the baseline RCI for each field. We calculated the fraction of each posterior distribution under each future scenario that were higher than or lower than a baseline range, defined as 95% to 105% of the field's historic median corn yield. We refer to these outcomes for each field-scenario as "upside" and "downside" probabilities. We then used downside probabilities to calculate absolute and relative risk mitigation scores in each weather scenario to compare two hypothetical management scenarios. Similarly, we used upside probabilities to calculate absolute and relative opportunity scores. We note that historical median field yield was used as a baseline for historical average corn yield for each field that is not heavily influenced by outliers.

Predictive Scenario Practice Name	1 Crop	2 Crop	3 Crop	4 Crop	5 Crop	6 Crop
Representative RCI Value	0	2.24	3.1	3.95	4.5	5.2

Table 2.2.. Names of rotational complexity scenarios by crop number with associated representative RCI value used for prediction.

More specifically, when comparing a simple rotational practice with one that is more diverse, we can use the downside probabilities to calculate the expected mitigation in risk offered by using the more diverse practice. Let d_s and d_D be the 95% downside probability for a chosen field in a given weather condition, under simple and diverse rotational management practices, respectively. Then the absolute risk mitigation offered by increasing from the simple rotation to the more diverse rotation for this field is $d_s - d_D$. We then define the relative risk mitigation as $(d_s - d_D)/d_s$. For example, Figure 2.3 shows posterior distributions for a hypothetical field under two competing rotations. In this figure, the orange curve represents a posterior predictive distribution under the simple rotation. This sample field has a median field yield of 20 Mg ha ⁻¹ which means 95% and 105% median yield for this field are 19 and 21 Mg ha ⁻¹, respectively. The figure shows downside probability of 0.2 under the simple rotation and 0.15 under the more diverse rotation, meaning the absolute risk mitigation would be (0.2 - 0.15)/(0.2 = 0.05)/(0.2

0.05/0.2 = 0.25 (a 25% relative reduction in risk). Under these definitions, positive values for risk mitigation represent situations where the more diverse rotation, D, has reduced risk compared to the simpler rotation, S.



Figure 2.3. Hypothetical predicted yield distributions for a single field under a simple rotation (orange, left) and a more diverse rotation (blue, right). This field has a historical average of 20 Mg ha⁻¹ making the cutoffs for 95% and 105% historical average yield 19 and 21 Mg ha⁻¹, respectively. The simple rotation has 95% downside probability of $d_s = 0.2$ whereas the diverse rotation has a smaller 95% downside probability $d_D = 0.15$, corresponding to an absolute risk mitigation score of $d_s - d_D = 0.2 - 0.15 - 0.05$ or a 5% absolute reduction in risk, and a relative risk mitigation score of $(d_s - d_D)/d_s = (0.2 - 0.15)/0.2 = 0.25$, or a 25% relative reduction in risk. Similarly, the simple rotation has a 105% upside probability of $u_s = 0.05$ whereas the diverse rotation has a larger 105% upside probability $u_D = 0.07$, corresponding to an absolute opportunity increase score of $u_D - u_s = 0.07 - 0.05 = 0.02$ or a 2% absolute increase in opportunity, and a relative opportunity increase score of $(u_D - u_s)/u_s = (0.07 - 0.05)/0.05 = 0.4$, or a 40% relative increase in opportunity.

Similarly, we define absolute and relative opportunity scores to compare the probability

of high crop yield under two different management practices. Let u_S and u_D be the 105% upside probability for a chosen field in a given weather condition, under simple and diverse rotational management practices, respectively. We then define absolute opportunity increase associated with using the more diverse rotation instead of the simpler rotation as $u_D - u_S$ and likewise the relative opportunity increase is $(u_D - u_S)/u_S$. For example, for the field represented in Figure 2.3, the 105% upside probabilities are $u_S = 0.05$ and $u_D = 0.07$. The associated absolute and relative opportunity increase scores would then be 0.07 - 0.05 = 0.02 (a 2% absolute increase in opportunity) and (0.07 - 0.05)/0.05 = 0.4 (a 40% relative increase in opportunity), respectively. In this case, positive values for opportunity increase represent situations where the more diverse rotation has increased probability of opportunity compared to the simpler rotation.

These risk mitigation and opportunity increase scores can be calculated to compare a variety of rotational practices over the three weather scenarios. In this way, we are able to quantify the field-level risk mitigation and opportunity increase afforded by using more diverse rotations over simpler rotations. This is particularly useful from the farmer and farm lending perspectives, as we can demonstrate the utility of using more diverse rotation in terms of risk mitigation and opportunity increase under various weather conditions and provide an economic rationale for increased rotational complexity.

2.2.6 Model Validation and Limitations

It was imperative to assess the performance of our model to ensure its robustness, model fit, and the accuracy of its predictions. In our early frequentist analysis, we compared models with a variety of predictors using measures such as R². When working in a Bayesian framework, a natural measure of predictive power for Bayesian models is empirical coverage probability (ECP). This is calculated by using the Bayesian model to create 95% credible intervals for insample predictions and then calculating the proportion of data points that lie within the credible intervals. We also used a leave-one-year-out validation approach to examine out-of-sample prediction by fitting models excluding a single year and then making predictions in that year.

To ensure sufficient modeling sample sizes, predictor variation, and scale of downside and upside probabilities, we imposed various restrictions on our modeling and prediction. First,

the minimum sample size for a county is 500 data points to avoid attempting to fit a model on counties with very few fields. The minimum number of fields modeled in a county is 143 and the median is 3013.

Second, since July maximum VPD is our primary weather predictor and we are predicting for values between 18 and 22 hPa, we decided to conservatively only model and predict for counties where the average July maximum VPD over the county exceeds 21 hPa in a single year (excluding 2012, an extremely dry year). In contrast to Illinois, Minnesota generally experiences less extremes and variability in VPD. Past research has shown that the detrimental effect of July maximum VPD on corn is approximately linear between 20 and 40 hPa whereas VPD levels below 20 hPa have minimal effect (Xu et al., 2021). This restriction allows modeling for all of Illinois and most of southern Minnesota while excluding counties in Minnesota where July maximum VPD of 20 and 22 hPa are uncharacteristic of the local conditions.

Finally, after calculating risk mitigation and opportunity scores for individual fields, when summarizing downside and upside posterior probabilities at the county level, we exclude fields from each county whenever the downside or upside probability predictions fall outside of [0.05, 0.95]. The justification for this choice of which fields to exclude is two-fold. First, downside and upside probabilities on the tails of the distribution (i.e., the bottom 5% and top 5% of the posterior predictive distribution) are less accurate (more prone to relative errors in estimating chances of events) than those in the bulk of the distribution. Further, these extreme probabilities are less meaningful for farmers and insurers when comparing practices and would generate misleading relative risk mitigation and opportunity scores. For example, suppose that in a given field for a given adverse weather scenario, our model predicts 99% downside probability

(probability of falling below 95% average field yield) under a simple rotation and 98% downside probability under a more diverse rotation. This would result in absolute risk mitigation of 1% and relative risk mitigation score of approximately 1%. In both these scenarios, the downside event is nearly certain to occur, and the management system becomes irrelevant. On the opportunity side, suppose that for this field, under an unfavorable weather scenario, the upside probability (probability of achieving above 105% average field yield) is 1% under a simple rotation and is 2% under a more diverse rotation. This would result in an absolute opportunity increase of 1% and relative opportunity score of 100%. However, the chance of achieving this upside is negligible in both scenarios, no farmer would expect to get that lucky, the difference between the management systems also becomes irrelevant, and reporting a 100% opportunity score would be highly misleading. The reader can imagine examples of scenarios for both upside and downside probabilities where the tails are reversed compared to the above two scenarios, and the conclusions of avoiding those fields in any reports and county-level summaries remains the same. Thus, to avoid all these extremes, we exclude these tail probabilities. We make a note, and we report the proportion of fields per county which are excluded by this restriction.

2.3 Results

2.3.1 Field-level risk mitigation and opportunity increase

Since our goal is to quantify risk mitigation for individual farmers, we begin with a fieldlevel presentation of results. We use Logan County to serve as an example. As discussed in Section 2.2.4, Bayesian Mixed Effects Model, we fit a model for the neighborhood around Logan County encompassing all counties adjacent to the focal county. The Bayesian coefficient estimates with associated 95% credible intervals are given in Table 2.3. Note that since all predictor variables are standardized, each coefficient estimate represents the change in yield associated with an increase of one standard deviation in the predictor variable. For example, we see that July maximum VPD has a large detrimental effect on yield as expected, with one standard deviation increase in July maximum VPD being associated with an estimated decrease in yield of 4.023 Mg ha⁻¹. We can use similar interpretations for each of the predictors to show general trends over the county. For example, we see that corn yield is estimated to increase both over time (year) and with better soil quality (NCCPI). Further, we see that field baseline RCI over the study period has a small estimated negative effect on yield but yearly deviation from baseline RCI has a positive estimated effect on yield of twice the magnitude. Further, coefficients for interaction terms between decomposed RCI terms and July maximum VPD are positive, meaning that higher RCI terms are associated with higher yield as July maximum VPD increases.

Parameter	Estimate (Mg ha ⁻¹)	Est. Error	95% Credible Interval
Intercept	21.816	0.0117	(21.794, 21.839)
Baseline RCI	-0.112	0.0128	(-0.137, -0.0873)
Yearly deviation from baseline RCI	0.226	0.0099	(0.207, 0.246)
NCCPI	0.419	0.0094	(0.400, 0.438)
July VPD Max	-4.028	0.0064	(-4.041, -4.015)
Year	0.437	0.0017	(0.433, 0.440)
Yearly deviation from baseline RCI x July VPD Max	0.167	0.0104	(0.146, 0.188)
Baseline RCI x July VPD Max	0.014	0.0084	(-0.0025, 0.003)
Baseline RCI x NCCPI	0.131	0.0106	(0.110, 0.152)
sigma	2.814	0.0047	(2.805, 2.823)

Table 2.3. Coefficient estimates, estimated error, and 95% credible intervals for model coefficients in Logan County.

After fitting our Bayesian model over the neighborhood, we then made predictions for each field under three weather conditions and seven rotational complexities resulting in 21 predicted scenarios and calculated 95% downside and 105% upside probabilities for each scenario. Table 2.4 shows predicted upside and downside probabilities for a single field in Logan County. For example, our model predicts that by using a rotation with RCI of 2.24 in normal conditions, a farmer has a 27.6% chance of falling below 95% median field yield and a 38.57% chance of achieving over 105% median field yield. We then see that as the number of crops in rotation increases, the 95% downside probability decreases and the 105% upside probability increases within each weather scenario. We also see that as July maximum VPD increases and we predict

RCI used for prediction	Number of Crops in Rotation (in a 6-year period)	Weather condition	95% Downside Probability	105% Upside Probability
0	1 Crop	Normal (July Max VPD 18 hPa)	0.2760	0.3857
2.24	2 Crops	Normal	0.2413	0.4347
3.1	3 Crops	Normal	0.2007	0.4840
3.95	4 Crops	Normal	0.1937	0.5053
4.5	5 Crops	Normal	0.1663	0.5333
5.2	6 Crops	Normal	0.1573	0.5423
	Historical Field Average (~2 Crops)	Normal	0.2343	0.4407
0	1 Crop	Somewhat Dry (July Max VPD 20 hPa)	0.5313	0.1663
2.24	2 Crops	Somewhat Dry	0.4207	0.2487
3.1	3 Crops	Somewhat Dry	0.3860	0.2823
3.95	4 Crops	Somewhat Dry	0.3510	0.3170
4.5	5 Crops	Somewhat Dry	0.3237	0.3410
5.2	6 Crops	Somewhat Dry	0.3073	0.3553
	Historical Field Average (~2 Crops)	Somewhat Dry	0.4197	0.2360
0	1 Crop	Dry (July max VPD 22 hPa)	0.7830	0.0467
2.24	2 Crops	Dry	0.6720	0.0927
3.1	3 Crops	Dry	0.6097	0.1257
3.95	4 Crops	Dry	0.5570	0.1563
4.5	5 Crops	Dry	0.5527	0.1673
5.2	6 Crops	Dry	0.4713	0.2127
	Historical Field Average (~2 Crops)	Dry	0.6437	0.1057

Table 2.4. 95% downside and 105% upside probabilities in each prediction scenario for a sample field in Logan County.

more dry scenarios, the 95% downside probability increases and the 105% upside probability decreases since there is a detrimental effect of high July maximum VPD on yield. Note that we also include predictions in each weather condition using the field average RCI to represent the downside and upside probabilities under the current management practice. For brevity, we classify this field average with a number of crops but note that the field average may differ slightly from the representative values used to predict for a certain number of crops.

From these predicted probabilities, we can calculate absolute and relative risk mitigation scores to compare two competing management conditions. Since 2-crop corn-soy rotations are most common in our study area, we use the scenario representing 2 crops as our simple rotation, S, and the scenario representing 3 crops as our more diverse rotation, D. Table 2.5 shows calculated risk mitigation scores for comparing these two rotational practices for our sample field.

Weather Condition	95% Downside Probability, d _s (RCI = 2.24, 2-crop)	95% Downside Probability, d _D (RCI = 3.1, 3-crop)	Absolute Risk Mitigation Score d _s - d _D
Normal	0.2413	0.2007	0.2413 - 0.2007 = 0.0406
Somewhat Dry	0.4207	0.3860	0.0347
Dry	0.6720	0.6097	0.0623

Table 2.5. Calculated absolute and relative risk mitigation scores for the sample field in Table 2.4 when comparing a 3-crop rotation over a 2-crop rotation.

We can then summarize and view such risk mitigation scores to understand the effect of using one practice instead of another over the entire county. Figure 2.4 shows box plots of the distribution of absolute and relative risk mitigation scores for Logan County when choosing a more diverse rotation (RCI of 3.1, 3 crops) instead of a simpler rotation (RCI of 2.24, 2 crops)

and Table 2.6 includes percentiles for the relative risk mitigation scores by field in each weather scenario. Results show that in all weather conditions, the 10th percentile of absolute risk mitigation scores is positive, meaning that in each condition, 90% of fields would benefit from using a 3-crop rotation over a 2-crop rotation. Further, in Figure 2.4 we can see visually that nearly all fields will benefit from the more complex rotation in this comparison, with increased absolute risk mitigation as weather becomes more dry.



Figure 2.4. Box plots showing the distribution of absolute (top) and relative (bottom) risk mitigation scores for all fields in Logan County. It is evident that as weather conditions become more dry, absolute risk mitigation increases, while relative risk mitigation stays approximately the same.

Percentile	10th	20th	30th	40th	50th	60th	70th	80th	90th	NA
Normal	0.0057	0.0097	0.0127	0.0153	0.018	0.021	0.024	0.028	0.034	223
Somewhat Dry	0.0193	0.0253	0.0297	0.0327	0.036	0.0393	0.043	0.0467	0.0527	41
Dry	0.031	0.0377	0.042	0.0457	0.0497	0.053	0.0567	0.0613	0.067	23

Table 2.6: Percentiles of absolute risk mitigation scores for Logan County in each weather scenario. Note that the number of fields excluded due to downside probabilities outside of [0.05, 0.95] is given under NA. For example, there are 223 fields without relative risk scores under normal conditions because in such conditions, many fields have a chance of falling below 95% median field yield that is below 5%. Similarly for 105% upside probabilities, we can create absolute and relative opportunity scores to

show the change in bumper crop opportunity when implementing one practice over another.

These results for our sample field are included in Table 2.7 and field-level results over Logan

County are summarized in Figure 2.5 to show the change in opportunity when choosing a 3-crop

over 2-crop rotation.

Weather Condition	105% Upside Probability, u _S (RCI = 2.24, 2-crop)	105% Upside Probability, u _D (RCI = 3.1, 3-crop)	Absolute Opportunity Increase Score u _D - u _S
Normal	0.4347	0.4840	0.0493
Somewhat Dry	0.2487	0.2823	0.0336
Dry	0.09267	0.1257	0.0330

Table 2.7. Calculated absolute and relative opportunity increase scores for the sample field in Table 2.4 when comparing a 3-crop rotation over a 2-crop rotation.



Figure 2.5. Box plots showing the distribution of absolute and relative opportunity scores for all fields in Logan County. It is evident that as weather conditions become more dry, relative change in opportunity goes up, while absolute change in opportunity stays approximately the same.

2.3.2 County-level summaries of risk mitigation and opportunity increase

To have a large-scale view of results over the entire study area, we used two methods to aggregate and display field-level results for risk mitigation and opportunity scores to the county level. The first method is to give the median absolute and relative risk mitigation score for each county. This represents the risk mitigation the "average" farmer would experience when comparing two competing practice levels. To connect back to section 2.3.1 on field-level results,

the median absolute risk mitigation score for Logan County is 0.0497 in dry conditions, according to Table 2.6, which can be seen visually as the center line of the dry boxplot in Figure 2.4. Figures 2.6 and 2.7 show the median absolute and relative risk mitigation, respectively, when using 3 crops instead of 2. In these figures, we see that in all weather conditions, all modeled counties in IL have positive median absolute and relative risk mitigation scores meaning that the "average" field will have reduced risk when using the more complex rotation. In contrast, median absolute and relative risk mitigation scores in MN are generally positive in normal conditions, nearly zero in somewhat dry conditions, and slightly negative in dry conditions. We note, however, that our dry scenario is one that occurs in IL approximately every 3-4 years whereas July maximum VPD as extreme as 22 hPa is much less frequent in MN (approximately once every ten years). When comparing the three weather scenarios, Figure 2.6 shows that median absolute risk mitigation tends to increase throughout IL as July maximum VPD increases. The relationship between median relative risk mitigation scores and July maximum VPD, visualized in Figure 2.7, seems to follow the opposite trend. However, this is partially due to the magnitude of downside probabilities in each condition. That is, drier conditions have larger downside probabilities, and therefore the same absolute reduction in risk has a smaller relative magnitude compared to the magnitude of the downside probability. Overall, Figures 2.6 and 2.7 show that using a more complex rotation including 3 crops instead of 2 has a risk mitigating effect throughout IL, particularly in more dry conditions. The story is less straightforward in the modeled counties of MN, but in the conditions that are more common in MN (normal and somewhat dry), there is some risk mitigation and no significant increase in risk afforded by using the more complex rotation.



Figure 2.6. Median absolute risk reduction scores for all fields in a county when using a 3-crop rotation instead of a 2-crop rotation for that field, in normal (left), somewhat dry (middle), and dry (right) conditions. This shows the absolute risk reduction the "average" field will experience.



Figure 2.7. Median absolute risk reduction scores for all fields in a county when using a 3-crop rotation instead of a 2-crop rotation for that field, in normal (left), somewhat dry (middle), and dry (right) conditions. This shows the relative risk reduction the "average" field will experience.

The second method of aggregation is to show the proportion of fields in a county that have absolute or relative risk mitigation above a certain threshold. Figures 2.8, 2.9, and 2.10 use this method with thresholds of 0 and 0.05, respectively. That is, Figure 2.8 shows the proportion of fields in each county that have any risk mitigation when using 3-crop instead of a 2-crop rotation. This is valuable because it shows what proportion of fields in each county will benefit from using the more complex rotation in each weather scenario. We can see from this figure that
nearly all counties in Illinois have over 90% of fields having predicted risk reduction when using the more complex rotation. Further, modeled counties of Minnesota have a high proportion of fields experiencing risk mitigation in normal conditions and over half of the included counties have over 50% of fields experiencing risk mitigation in somewhat dry conditions.



Figure 2.8. The proportion of fields per county that have positive absolute and relative risk mitigation scores when using a 3-crop rotation instead of a 2-crop rotation. Note that by definition, a positive absolute risk mitigation score implies a positive relative risk mitigation score.

In contrast, Figures 2.9 and 2.10 use a threshold of 5% absolute or relative risk mitigation, respectively. This is valuable from an insurance perspective, as we can see what proportion of fields in a county will have risk mitigation that "moves the needle". We can use visualizations like Figures 2.9 and 2.10 to identify counties in our study area that would most benefit from widespread adoption of a more complex rotation. For example, it is clear from Figure 2.10 that targeting central would likely be the most profitable for a farm lender or insurer when encouraging widespread adoption of a 3-crop rotation over the common corn-soy rotation.



Figure 2.9. The proportion of fields per county that have absolute risk mitigation scores greater than 5% when using a 3-crop rotation instead of a 2-crop rotation.



Figure 2.10. The proportion of fields per county that have relative risk mitigation scores greater than 5% when using a 3-crop rotation instead of a 2-crop rotation.

Similarly, we can perform the same aggregations for absolute and relative opportunity increase, as shown in Figures 2.11, 2.12, and 2.13. Similar patterns are identified when evaluating aggregated visualizations of opportunity increase over the study area. For example, we see that throughout Illinois, the use of a more complex rotation is associated with positive absolute and relative opportunity increase scores, particularly in more dry conditions.



Figure 2.11. Median absolute opportunity scores for all fields in a county when using a 3-crop rotation instead of a 2-crop rotation for that field, in normal (left), somewhat dry (middle), and dry (right) conditions. This shows the absolute change in opportunity the "average" field will experience.



Figure 2.12. Median relative opportunity scores for all fields in a county when using a 3-crop rotation instead of a 2-crop rotation for that field, in normal (left), somewhat dry (middle), and dry (right) conditions. This shows relative change in opportunity the "average" field will experience.



Figure 2.13. The proportion of fields per county that have positive absolute and relative opportunity increase scores when using a 3-crop rotation instead of a 2-crop rotation. Note that by definition, a positive absolute opportunity score implies a positive relative opportunity score.

2.3.3 Presentation and interpretation of model coefficient estimates by county

Beyond quantifying the risk mitigation and opportunity increase under varying predicted scenarios, we can use the Bayesian estimates for various model coefficients to analyze the effects of rotational complexity, water stress, and soil quality in a regional way. For example, we can confirm expected effects of time and water stress over the entire study area. Figure 2.14 shows Bayesian coefficient estimates for July maximum VPD and year by county. We see that July maximum VPD has a strong estimated negative effect on yield, with more detrimental effects occurring in the Southern part of the study area. This estimated effect is weaker in Minnesota counties and may reflect the non-linear relationship between July maximum VPD and yield, since the detrimental effect of July maximum VPD on corn is linear between approximately 20 and 40 hPa and such extreme levels of July maximum VPD are much less common in Minnesota than Illinois. Figure 2.14 also shows that corn yield is estimated to increase by approximately 2 - 6 Mg ha⁻¹ per year, with a larger increase in Western Illinois and the Northern part of the studied counties in Minnesota.



Figure 2.14. Heat maps of Bayesian coefficient estimates in Mg ha⁻¹ for July Maximum VPD (left) and year (right) by county for the study area. July maximum VPD coefficient estimates show a large detrimental effect of high July maximum VPD on yield that increases in more Southern counties. Coefficient estimates for year show a small increase in yield over time, with greater increases in Western Illinois and the Northern portion of the studied area in Minnesota.

With the detrimental effect of July maximum VPD in mind, we can examine the estimated effects of increased rotational complexity and its interaction with water stress. Figure 2.15 shows coefficient estimates for baseline RCI and yearly deviation from baseline RCI, as well as their interactions with July maximum VPD. First examining coefficient estimates for baseline RCI, we see a generally small estimated negative effect associated with fields with higher baseline rotational complexity. This relationship may be confounded with the fact that historically, farmers on marginal lands generally employ higher rotations, and thus soil quality is also part of the explanation. Farmers may try to mitigate the negative effect on yield associated with lower quality soil by employing regenerative soil practices, resulting in a confounded relationship between baseline RCI, soil quality (NCCPI), and yield. Turning our attention to yearly deviations from baseline RCI, we see a generally positive estimated effect, meaning that

within a given farm, increasing rotational complexity compared to their usual practice is associated with a yield benefit. Looking at interactions between RCI terms and July maximum VPD, the regional trends are less universal. For the interaction between baseline RCI and July maximum VPD, when comparing fields under the same VPD conditions, positive values mean fields with higher baseline RCI will have mitigated risk in terms of yield losses due to water stress. As July maximum VPD increases, this reduction in risk due to VPD also increases. In Figure 2.15, we see regionally that Northern and Central Illinois as well as modeled counties in Minnesota experience these positive values, meaning that fields with higher baseline RCI are predicted to experience greater risk mitigation in periods of water stress than comparable fields with lower baseline RCI. Finally for the interaction between yearly deviation from baseline RCI and July maximum VPD, when comparing fields under the same July maximum VPD conditions, positive values mean fields with larger RCI increases year-to-year are predicted to experience mitigated risk in terms of yield losses due to July maximum VPD, with this risk reduction increasing as July maximum VPD increases. In Figure 2.15, we see regionally that Northern and Central Illinois experience these positive values, representing mitigated risk due to dry weather when farmers increase from their historical rotational complexity.

Baseline RCI

Yearly Deviation from Baseline RCI





-0.2 to 0.6

0.4

0.0



Interaction between July Max VPD and Baseline RCI

Interaction between July Max VPD and Yearly Deviation from Baseline RCI



Figure 2.15. Heat maps of coefficient estimates for baseline RCI (top left) and yearly deviation from baseline RCI (top right), as well as their interaction with July Maximum VPD (bottom left and right, respectively).

Finally, we can view the regional effect of soil quality (NCCPI) and its interaction with baseline RCI in Figure 2.16. As expected, higher soil quality is associated with increased yield.

In Central to Northern IL and in MN when comparing fields with the same soil quality, fields with a higher baseline RCI will have a small estimated increase in yield, corresponding to positive coefficient estimates.



Figure 2.16. Heat maps of coefficients for NCCPI (left) and the interaction between NCCPI and baseline RCI (right). As expected, higher soil quality is associated with increased yield. In Central to Northern Illinois and in Minnesota when comparing fields with the same soil quality, fields with a higher baseline RCI will have an associated small increase in yield.

We can also compare coefficient estimates by state through a box plot of coefficient estimates by county in Figure 2.17. From the figure, we see that year and NCCPI have positive coefficients as expected, with yearly increases in yield estimated to be larger in Minnesota, and estimated positive effects of soil quality being generally larger in Illinois. The interaction between baseline RCI and NCCPI by county is not consistent regionally in Illinois, as the boxplot spans both sides of zero, but has a generally positive estimated effect in Minnesota. Confirming the complex story of interactions between July Maximum VPD and RCI presented in Figure 2.15, the box plots for these interaction terms in Illinois are on either side of zero, whereas the coefficients for these interactions in Minnesota lean positive for the interaction with baseline RCI and lean negative for the interaction with yearly deviation from baseline RCI. Finally, baseline RCI seems to have a generally negative estimated effect while yearly deviations from baseline RCI seem to have a generally positive effect, matching the trends from Figure 2.15. These relationships are more clearly to one direction in Illinois than in Minnesota. By examining the distribution of coefficient estimates by county separately for the two states, we are able to see visually how relationships between main predictor variables and yield differ between the two states.



Figure 2.17. Boxplots of coefficient estimates for all model predictors by county, separated by state. Here each boxplot represents the distribution of coefficient estimates for the given variable over all modeled counties in the corresponding state.

2.3.4 Model Validation

To evaluate the accuracy of model predictions, we calculate empirical coverage probability for prediction in each county using observed data for predictors. Figure 2.18 shows

the empirical coverage probability by county for 95% credible intervals. This is extremely accurate uncertainty quantification, as all coverage probabilities are near the nominal level of 0.95. In this figure, there is almost no under-reporting (i.e., coverage probability less than 95%) so we can feel confident we are not giving a false sense of accuracy. In counties with above the nominal level, we err on the side of conservative estimates. With these results, we can be confident in the accuracy of our uncertainty quantification.

Empirical Coverage Probability



Figure 2.18. The empirical coverage probability for prediction in a single county. That is, the proportion of actual yield data observations that fall within 95% credible intervals for yield created through model prediction using observed data for predictors.

2.4 Limitations and Future Work

2.4.1 Regional Variation

Though our analysis shows clear trends of risk mitigation and opportunity increase associated with increased rotational complexity throughout Illinois, the relationships between rotational complexity, weather conditions, and corn yield are less clear in Minnesota. A contributing factor may be the regional differences in July VPD, as Minnesota does not experience extremes above 20 hPa that are detrimental to corn as frequently. These differences in July maximum VPD also excluded a large proportion of Minnesota counties from our analysis. Further, planting dates vary regionally meaning that the stages of development where corn is most sensitive to water stress may not align in different regions. In future model iterations, we plan to explore the effects of VPD later in the year to try to account for differences in planting dates and to expand our analysis to other weather predictors that may be more appropriate in other regions. Another contributing factor to differences in results between the two states may be the difference in number of years of data availability. While our dataset encompasses yield observations in Illinois from 2005 to 2020, our Minnesota data is limited to 2011 to 2020. As our datasets expand to cover more field-years, we expect the accuracy of our modeling will continue to improve.

2.4.2 Model Expansion

Currently our work is limited to examining the effect of increased rotational complexity and dry conditions on corn yield. In future analysis, we plan to incorporate a variety of regenerative soil practices including conservation tillage and cover-cropping, as well as examining relationships between concurrent practices. Further, we will expand our analysis to

include other weather conditions such as flooding and add soybeans as an additional crop. We are also currently expanding our dataset to encompass nine states in the Midwestern US. We currently have yield data for both corn and soy and a variety of weather and soil variables as well as county-level statistics on management inputs and indemnity payouts. In this way, we can further quantify the risk mitigation associated with regenerative management practices under a variety of weather conditions and create the empirical link between these practices and reduced risk.

We currently examine 95% downside and 105% upside probabilities in 21 scenarios but want to use the full advantage of Bayesian predictive distributions to answer a variety of questions and compare other possible management options and weather conditions. We are currently limited by computational considerations in terms of time taken to model fit and make predictions and data storage capacity. With improvements in data storage, we plan to store full posterior predictive distributions in order to answer questions about a variety of possible outcomes (e.g., probability of dropping below 80% average field yield, etc.).

We want this work to be directly beneficial to farmers by supplying lenders with risk reduction metrics to translate the economic benefit of regenerative practices to reduced rates for farmers. To facilitate this, we are building an interactive tool for use by lenders to compare management practices and evaluate risk reduction based on our model results. Further, we are beginning to incorporate economic factors related to crop pricing to give more accurate assessments of the economic benefit of regenerative practices. We hope that by providing the economic rationale for adopting regenerative soil practices, we can help encourage widespread adoption of these soil-protecting measures.

2.5 Conclusion

As climate change increases the frequency and severity of adverse weather conditions, it is vital to implement farm management practices that can help prevent crop loss. Increased rotational complexity has been shown in case-study experiments to increase crop yield over time in average conditions and to mitigate the detrimental effects of harsh weather conditions like drought. In our analysis, we targeted important corn producers in the US corn belt and quantified the risk mitigation effect of adopting more diverse rotations at the field level, particularly in dry conditions, on a regional scale.

By using a Bayesian framework, we incorporated previous domain research on the effects of rotational complexity, soil quality, water stress, and time on corn yield. With our neighborhood modeling approach using a mixed effects model, we were able to account for regional variability in practice adoption, weather trends, and soil quality, as well as field-level differences in rotational complexity and overall productivity. Our unique methodology allows us to make comparisons between rotational practices at the field level and use past field history and field-level characteristics such as soil quality to make accurate predictions.

Our results show promising risk mitigation associated with higher rotational complexity, particularly in dry weather in Illinois. By performing field-level risk analysis, we can provide individual farmers and loan officers the information necessary to make informed decisions on practice adoption. Further, aggregated county-level risk summaries are valuable for lenders to prioritize practice adoption in areas with greatest risk reduction. Finally, visualization of coefficient estimates over the study area can help to identify key relationships between predictors and regional trends in weather variability, practice adoption, and soil quality.

Through this work, we have established an empirical connection between diverse crop

rotation and risk mitigation in dry weather on a regional scale. As our work continues, we hope to expand our modeling efforts to incorporate a larger geographical area as well as a variety of weather conditions, management practices, and soil characteristics. By expanding our analysis, we will continue to provide empirical evidence and economic rationale to support the widespread adoption of regenerative soil practices throughout the Midwestern US and beyond.

CHAPTER 3: "YOU CAN LEARN R": AN ACCESSIBLE AND INCLUSIVE WORKSHOP TO TEACH RESEARCH PROFESSIONALS HOW TO LEARN R

3.1 Introduction and Workshop Motivation

The statistical programming language R is widely used in research across disciplines by academics, students, and industry professionals worldwide (Worsley, 2022). One advantage of R is its free and open-source nature. Since its inception in the mid-90s, R has grown exponentially, boasting nearly 20,000 packages and widespread usage across academia and industry globally (R Core Team, 2021). However, despite its popularity, learning R can be challenging, with a steep learning curve that makes it inaccessible for many students (Gallagher, 2022). With the R language's rapid growth, numerous learning resources now exist for R, ranging from comprehensive online courses and textbooks to concise blog posts and videos. However, the sheer abundance of packages and resources can overwhelm researchers who are eager to learn R but struggle to identify the most suitable packages and learning materials for their research needs. This challenge is particularly prevalent among early-career researchers like graduate students and postdoctoral researchers, who lack the time to navigate the extensive array of resources or engage in lengthy courses or books on R. Consequently, self-directed, ad-hoc learning becomes the norm, leading to uncertain learning outcomes and potentially discouraging learners from pursuing R further (Theobold & Hancock, 2019). There is currently a need for a resource aimed at researchers who want to utilize the advantages afforded by learning R but lack direction in how to start their learning journey and which resources will be most beneficial for their specific learning and research needs.

To address this growing need, I developed an inclusive and accessible workshop aimed at teaching early-career researchers how to learn R. The workshop includes a seminar session, a curated resource document, and a working session with R exercises. Its goal is to provide

researchers with a starting point in learning R, boosting their interest, confidence, and ability to learn R for use in their research, and providing guidance in selecting appropriate learning resources tailored to the specific needs of each researcher. The workshop was successfully conducted with graduate students and post-docs at Michigan State University and served as a foundational module for a novel training experience on Bayesian methods in agronomy for earlycareer professionals in agriculture science in Africa. Surveys were used to gather feedback for further improvements and measure the workshop's success in enhancing participants' interest, confidence, and ability in learning R. This chapter will present the workshop materials, discuss the process and considerations involved in creating this inclusive and accessible resource, analyze the results of piloting the workshop, discuss the use of the workshop as a foundational module for a larger five-day workshop on Bayesian methods in agronomy, and outline future plans for improvements and future iterations of both workshops.

3.2 "You Can Learn R"

The "You Can Learn R" workshop is a multi-faceted learning experience including an inperson seminar with cooperative learning exercises, an online-hosted written document with advice for learning R and curated R-learning resources, and a working session to implement the learning from the seminar and written portions. The workshop was designed with accessibility and inclusion in mind with the primary aim to increase participant interest, confidence, and ability in learning R.

3.2.1 The seminar

The "You Can Learn R" workshop begins with an in-person seminar session to introduce participants to the R language and facilitate cooperative learning exercises to get participants

started working with R. The beginning of the seminar portion introduces the workshop purpose, funding, structure, and presenter. The body of the seminar includes three sections: "R: What, why, and how?", "Learning new techniques in R", and "When things go wrong". Each section includes a short lecture portion followed by practical exercises to encourage active and cooperative learning.

The first section introduces R and RStudio and motivates how a researcher might benefit from using R in their research. This section outlines the benefits of R such as its free and opensource nature, the ability to accomplish nearly any data-driven task with the vast library of packages, the welcoming and widespread R learning community, and the advantages of using R in terms of reproducible research. The section also explains how R and RStudio work together and how R consists of the built-in packages in base R, additional packages, and functions within those packages. This section concludes with an exercise where participants access and explore RStudio and execute a short R script to learn about running basic commands.

The second section gives advice for learning new techniques in R based on my experience teaching and learning R as well as recommendations from various R learning resources. The first piece of advice is motivating R work with research-related projects, sample project and data, or R challenges, and writing out the required steps explicitly. After finding motivation and planning what to accomplish in R, we recommend strategic searching practices to facilitate accomplishing the required steps for the motivating project. Strategic searching includes utilizing R help menus and package documentation, searching online using package names and specific sites as keywords, and copying and modifying existing R code examples to accomplish desired tasks. Further, we encourage participants to learn with others and find an R learning community that fits their specific needs. Finally, we recommend managing expectations

when learning R stating, "R can do anything, but you don't need to know it all."

After outlining this advice, participants are asked to use these recommendations to complete two exercises. Each exercise has two parts: part (a) provides sample code that performs a specific task and asks participants to use strategic searching tips to find out what each part of the code does and write out the steps in the code comments while part (b) provides steps to complete a related task and asks participants to adjust the code from part (a) to complete said task. In this way, participants are able to practice the recommendations given in the preceding lecture portion to learn what a coding example does and adapt that example to solve a motivating problem. The in-person and collaborative nature of the seminar encourages learning with others and gives participants a starting point for building their own R learning network. The two exercises use two common packages from the tidyverse, "an opinionated collection of R packages designed for data science [where] all packages share an underlying design philosophy, grammar, and data structures" (Wickham et al., 2019). The first exercise involves data manipulation with the package dplyr (Wickham et al., 2020) while the second exercise uses ggplot2 (Wickham, 2016) to perform data visualization. These packages were chosen specifically to introduce participants to two extremely useful and powerful packages for manipulating, summarizing, and visualizing data in R.

The final section of the seminar is titled "When Things Go Wrong" and presents common errors in R and how to troubleshoot when errors occur. This section is titled as such to convey to participants that errors are common when using the R language and they should not be discouraged when they inevitably make a mistake. Common errors are discussed including errors related to capitalization, misspelling, closing or continuing punctuation, conflicting code, unloaded libraries, and unsaved objects. Participants are then given advice for troubleshooting

such as strategic searching, running code line-by-line, and asking other R users for help. Participants then engage with the final set of exercises for the seminar. In these exercises, users are presented with numerous, nearly identical chunks of code, each with a single change that will result in an error. This section again includes two exercises, each based on the dplyr and ggplot2 code chunks presented in the previous section's exercises. Both exercises include a final challenge code chunk where participants are encouraged to create errors of their own for the presenter to troubleshoot in front of the group. In this way, participants can see in real time how a more seasoned R user works through error messages.

After the conclusion of the seminar, participants are given access to the "You Can Learn R" written document and encouraged to return for a group working session to engage with the written document and collaborate with other R learners to accomplish tasks related to their own research.

3.2.2 The written document

The R language cannot be taught within a single seminar and the "You Can Learn R" seminar portion is merely a starting point in each participant's journey in learning R. To guide participants in their R learning, I created a written document with advice and curated resources for learning R. This resource was created in bookdown (Xie, 2020) and is hosted at https://www.bookdown.org/manskisa/You Can Learn R. The written document opens with a preface motivating the creation of the resource, outlining the target audience, and expressing the caveat that the document will continue to grow and expand over time. Chapter 1 relays the content of the seminar portion of "You Can Learn R" and provides a link to a cloud-hosted version of the seminar R project. Chapter 2 is a curated list of R resources divided into three sections: recommended packages, learning resources, and learning communities. The first

section gives a list of commonly used packages with descriptions to point users toward potential helpful tools for their research. The second section includes recommended learning resources including written materials, interactive tutorials, data sources, and sites to search. These resources only encompass a small number of the expanse of learning materials related to R. However, they offer a few places to get started to help learners avoid the overwhelm of too many options with nowhere to start. The third section gives recommended learning communities. The worldwide R learning community is extensive and welcoming to learners at all levels. This section includes global communities such as R-ladies and the R for Data Science community, as well as communities for under-represented R user groups such as AfricaR, R-ladies, and Minorities in R. Chapter 2 ends with an appendix pointing to numerous R package cheatsheets. While cheatsheets are not recommended for learning R, they are designed for aiding quick understanding and use of functions and can be useful as a quick reference when working in R. The final chapter of the "You Can Learn R" written document discusses accessibility recommendations. This currently includes sections with tips for learning R with limited internet access, learning materials that have been translated to various languages, and recommendations for blind R users.

3.2.3 The working session

To allow participants time to work with the written document in a self-directed way with access to other R users, we offer a later working session where participants can use the provided resources alongside other learners to solve their own data-driven tasks. These tasks are openended and could include completing exercises based on the resources in the written document, following along with provided sample coding materials, or working on tasks related to the participants' own research. This session allows participants to put into practice the advice given

in the seminar and written document in a cooperative learning environment and apply that knowledge to data-driven tasks specific to their needs.

3.3 Workshop development considerations and peer feedback

This workshop represents a "Teaching as Research" project with the research question: "How does an accessibility and inclusion-based workshop for learning R affect researching professionals' interest, confidence, and ability to learn and use new techniques in R?" Thus, the "You Can Learn R" workshop aims to be an accessible and inclusive resource to teach strategies for how to learn R for a wide audience of researchers across disciplines. That is, the workshop aims to eliminate potential barriers to R learning based on ability and engage and include researchers from a variety of backgrounds. To facilitate this goal, numerous considerations were made in the development of this workshop to ensure continued accessibility and inclusion. Specifically, we follow the two broad goals outlined by Dogucu et al. (2023) in their framework for accessible and inclusive teaching materials for statistics and data science courses: "Goal 1. Course materials should be physically accessible" and "Goal 2. The development and delivery of course materials should be inclusive of a diverse body of learners" (p. 2). This section outlines design considerations that were made to answer the research question while following the above goals.

3.3.1 Development considerations

Teaching researchers HOW to learn R. A distinction that separates "You Can Learn R" from other R learning resources is that this workshop is not designed to teach participants R, but to teach them how to learn R. Each R user's learning experience is different and my experience has shown that many researchers learn R in a self-directed way. "You Can Learn R" does not claim to teach participants R but rather offers advice and direction on how participants can learn R themselves. This philosophy follows the overarching research goal of increasing participant interest, confidence, and ability to learn and use new techniques in R.

A multi-faceted workshop. The workshop design includes three portions to maximize the lasting impact on participants and allow flexibility in the learning process. The first portion is an in-person seminar where participants are introduced to R and engage in cooperative learning exercises with other participants. After the seminar, participants are given access to a written document of curated R learning resources and recommendations to guide them as they continue learning R. Finally, participants are invited to an in-person working session where they can use the knowledge and materials from the seminar and written document to tackle exercises related to their own research. In this way, the "You Can Learn R" workshop not only introduces participants to R and helps them get started in learning, but also offers a guide for their continued learning and a supervised and cooperative working opportunity to attempt using R in their own work.

An R workshop made IN R. To demonstrate the utility and flexibility of R, the "You Can Learn R" workshop was fully developed in R. All resources for the seminar portion are contained within an R project accessible in Posit Cloud. The seminar uses rmarkdown for the slide presentation, exercises, and solutions (Allaire et al., 2023; Xie et al., 2018, 2020). The R written document was created in bookdown, an R resource for creating and publishing books in R (Xie, 2020). Furthermore, participant survey results were analyzed, summarized, and visualized using R.

Human subject research. To measure the success of the workshop in increasing participant interest, confidence, and ability learning R, it was vital to be able to survey participants

throughout the workshop and analyze changes in their interest, confidence, and ability learning R over time. Since this work is research with human subjects, a project proposal was submitted to the MSU Internal Review Board for review and the research was determined exempt. All materials including study design, recruitment emails, and surveys were submitted for review and approved.

An open-access workshop. All workshop materials have been developed such that they are openly and continually accessible to participants. The seminar portion is hosted in Posit Cloud and the written document is hosted in bookdown, both freely accessible and downloadable online.

A living workshop. The workshop materials for "You Can Learn R" are intended to grow and change over time. With the vast number of R learning resources and the breadth of R packages growing every day, a static workshop cannot be expected to be sufficient as time passes. Further, one workshop developer cannot be expected to know and include all the valuable R learning resources and packages. Therefore, this workshop is intended to grow and change as more iterations are run based on participant feedback. For the written document, participants are encouraged to submit to the author recommendations for resources or sections they would like added.

RStudio and Posit Cloud. This workshop focuses on R using RStudio and its companion online version, Posit Cloud. These development environments are openly accessible, with the desktop version of RStudio being available for download on Windows, Mac, and Linux machines. Posit Cloud offers an online alternative to RStudio where projects can be easily shared between users. All seminar materials are hosted in Posit Cloud so participants can access the presentation, exercises, and exercise solutions at any time.

Tailored to a broad audience. This workshop is intended to be accessible for a wide variety of researchers that could benefit from using R in their research, but may not know where to start or have the time for extensive courses on learning R. To address these constraints, this workshop was designed to require minimal time commitment from participants and offer resources that can point participants in the best direction for learning the R skills necessary to complete tasks related to their research. Further, the written document addresses specific learning needs such as how to learn R with limited internet access and which learning resources are translated into non-English languages.

Inclusion through cooperative learning exercises. Under Goal 2, Dogucu et al. (2023) recommend the following strategies:

- "Showcase the diversity of the field through a broad group of scholars
- Use inclusive language, assumptions, and examples
- Use active learning approaches that encourage students to learn by doing
- Embrace the challenges and failures which are critical to learning
- Build rapport"

By holding the seminar portion of the workshop in-person, we brought together a group of scholars from a variety of disciplines and fostered an active-learning environment where participants learn by doing. Specific exercises focused on learning new techniques in R and error-handling, giving participants an opportunity to embrace the challenges and failures commonly encountered when learning R in a safe and supportive environment with other R learners. Since the first full run of the workshop was for a diverse group of early-career agronomy professionals in Africa, exercises were designed to be inclusive of the specific audience. These exercises used sample agronomy data and involved data manipulation and

visualization of African livestock data to demonstrate the utility of R in agronomy and commonly used techniques for researchers working with data. Finally, the in-person nature of both the seminar and the working session allowed time for participants and instructors to build rapport and led to participants being comfortable and eager to ask questions and work collaboratively.

3.3.2 Soliciting peer feedback

To evaluate the adequacy of the workshop at meeting the above goals in terms of accessibility and inclusion, a development version of the "You Can Learn R" seminar was presented to solicit peer feedback from members of the following groups at MSU:

- graduate students and faculty in the department of statistics and probability, for their feedback as statistics educators and experienced R users;
- fellows from the Future Academic Scholars in Teaching (FAST) fellowship, for their feedback as educators and novice R users;
- fellows from the Great IDEA fellowship, for their feedback on workshop accessibility and inclusivity; and
- graduate students from the Graduate Student Accessibility and Support Network (GSASN), for feedback on workshop accessibility.

These peer reviewers were each given feedback forms that offered some background on the workshop and questions related to workshop content, teaching and design, and surveys. Running a preview session was essential to ensuring the workshop was accessible and inclusive for the desired audience: early career researchers such as graduate students and postdoctoral scholars from a variety of disciplines with limited or no R experience. For example, based on feedback from peer reviewers, it was evident that early versions of the exercises were too difficult for

beginning R users. Using thoughtful feedback from these peer groups, we could improve the accessibility and inclusivity of the resource before piloting on the target audience.

3.4 Presentation and evaluation of results

3.4.1 Surveys

To measure the success of the workshop in terms of increasing participant interest, confidence, and ability in learning R, participants completed surveys at three points during the workshop: before the seminar portion, after the seminar, and after using the written document in a working session to complete tasks in R. These surveys included statements on a 7-point Likert scale related to interest, confidence, and ability in R, with scale response options including Strongly disagree, Disagree, Somewhat Disagree, Neutral, Somewhat Agree, Agree, and Strongly Agree. These statements were based on various technology usability and user experience surveys including USE (Lund, 2001), the Unified Theory on Acceptance and Use of Technology (UTAUT) (Venkatesh et al., 2003), and the System Usability Scale (SUS) (Brooke, 1995). Each survey had blocks of 3 to 8 related Likert questions according to the groupings of the established user experience surveys. These groupings include user intent to use R, usefulness of R, and satisfaction using R to measure participant interest in the R language. Groups to measure participant confidence in using and learning R include ease of and difficulty using R, ease of learning R, learning needs, and perceived ability to complete tasks. These statements remained consistent over the three surveys to facilitate comparisons across the different time points. The first survey included additional questions on participants' previous experience with R; other statistical tools such as Microsoft Excel, Stata, and SAS; and statistics. Surveys 2 and 3 included additional questions to solicit participant feedback on the workshop in order to inform

future improvements and evaluate the accessibility and inclusivity of the resource.

3.4.2 Presentations

A pilot presentation of the "You Can Learn R" seminar was given to graduate students and postdoctoral researchers from a variety of disciplines at MSU. The seminar was scheduled for 90 minutes and consisted of a 20-minute introduction where participants completed the first survey, three 20-minute blocks for each of the three seminar sections, and 10 minutes at the end for the second survey. The seminar was held in-person in a 32-seat computer lab where each participant had a computer provided. Despite nearly 100 researchers registering to attend the workshop, the capacity limitation meant that only 32 interested individuals were invited. Of the 32 invited, 13 individuals attended and completed the first survey, ten people completed the second survey, and only three people signed up to attend the follow-up working session. Due to this attrition and time limitation, the working session was not held in this pilot iteration of the workshop.

The workshop was first presented in its entirety to early-career researchers in agronomy as the foundational module for a 5-day learning experience on the fundamentals of Bayesian statistics in agronomy, offered in Addis Ababa, Ethiopia. The workshop was allotted an 80minute session in the morning followed by two approximately 2-hour sessions in the afternoon. The intent was to split the seminar portion between the first two sessions and use the final afternoon session as the working session. However, participants were so engaged when working on the various exercises with their peers that we extended each section of the seminar to comprise one of the three allotted sessions. We finished this day by distributing the second survey and the written document. In total, we had over 30 participants with 29 completing the first survey and 27 completing the second survey. Based on feedback from the second survey, we

found that some participants that were absolute beginners in R still struggled completing some of the exercises. To offer additional assistance, we held a 2-hour office-hour session where participants could work together in R while having access to the presenter to ask questions. Later in the week, we held a working session where participants formed groups, followed sample code to perform Bayesian analysis on their own data, and presented their results to the larger group.

3.4.3 Participant demographics and previous R exposure

Participants from both MSU and Addis Ababa consisted of early-career research professionals. At MSU, the participants represented a diverse range of disciplines, including mathematics; communicative sciences and disorders; human resources and labor relations; microbiology and molecular genetics; agricultural, food, and resource economics; plant, soil, and microbial science; plant pathology; chemistry; cell and molecular biology; and supply chain management. In Addis Ababa, workshop attendees were part of a larger workshop on Bayesian statistics in Agronomy, which was advertised throughout the CIMMYT Excellence in Agronomy (EiA) initiative network and National Agricultural Research System (NARS) partners across Africa, with a particular emphasis on encouraging junior scientists to apply. Attendees included agronomy researchers from various agricultural research organizations in Ethiopia and across Africa (see Section 3.5 for specific organizations).

The participants who completed Survey 1 (42 in total) had varying levels of prior experience with R, which could be classified into two categories: little to no experience (43%) and self-directed or contextual learning for specific purposes (57%). Those categorized as having little to no experience either responded "no" when asked if they had prior experience using R or claimed to have very little, some, or basic experience. The remaining participants had some experience with R, but their experiences aligned with previous observations on how early-career

researchers typically use and learn R. Learning was often self-directed and aimed at accomplishing tasks within their own research or studies. Many participants mentioned using R for general purposes such as data analysis, manipulation, visualization, and spatial analysis. Some provided specific examples or techniques related to their respective fields, such as "for agronomic soil properties data analysis and display," "to computationally analyze the data from flow cytometry," "to analyze some RNA sequence data," and "for linear programming." A few participants mentioned exposure to R during their studies, stating that they used R for thesis work, graduate study, or in a past course. Some participants also mentioned limited formal R training, with varying levels of learning outcomes, such as attending a basic crash course and acquiring a few coding skills or completing entry-level datacamp courses on using R for basic plotting and data manipulation. However, uncertainty in learning outcomes and ad-hoc selflearning were common themes observed among researchers. For instance, one participant stated, "I have used R to perform principal component analysis for my research. I have also attempted PC regression and partial least squares with a lot worse results," while another mentioned, "I have used tidyverse/dplyr/ggplot2 to do basic data cleaning/manipulation and to make plots; however, I'm not sure that my scripts are always efficient or follow 'best practices' in these areas." Some participants described their prior experiences with R as learning from online resources, such as YouTube tutorials or graphing through self-directed learning over a year. When asked to give examples where they could successfully complete tasks in R, responses such as the following spoke to the types of resources participants used to learn R before the seminar:

• "For data visualization and analysis through training received from colleagues and internet I became improved in using ggplot2, tapply, multicompview, among other packages"

- "There are a wide range of resources and forums to call upon for support/ assistance,
 i.e. R for Data Science, Stack Overflow, etc"
- "Stack Overflow does an amazing job"

These resources align with both recommended techniques (strategic internet searching and asking colleagues for help) and recommended resources (R for Data Science, Stack Overflow) that are included in the "You Can Learn R" workshop. Based on the participants' responses regarding their experiences with R, it is evident that these researchers fall within our target audience and have encountered similar learning experiences with R as those observed during my time teaching R.

3.4.4 Evaluating results and workshop feedback

To evaluate the success of the "You Can Learn R" workshop, quantitative and qualitative survey results related to participant interest, confidence, and ability in learning R were analyzed. Further, qualitative feedback from participants on the most valuable aspects of the seminar and where to improve, as well as on the accessibility and inclusivity of the workshop, helped to inform adjustments to the workshop and to evaluate the appropriateness of the included materials.

In terms of interest, confidence in learning R, we can examine the distribution of responses to Likert statements related to interest and confidence between Survey 1 and Survey 2 to evaluate the changes in participant attitudes between before and after the seminar portion. For interest, Likert statements were grouped into three areas: intent to use R, usefulness of R, and satisfaction using R. The distribution of responses to statements in these groups are summarized in Figure 3.1. From the figure, we see that statements on intent to use R and satisfaction using R had a higher proportion of participants in agreement after the seminar portion (Survey 2) than

before (Survey 1). The change in perceived usefulness of R between the two surveys is less clear, and may be attributed to the fact that attending the seminar gives participants an understanding of both the utility and complexity of R. For example, when asked how the seminar affected participant interest in R, some expressed a positive change in interest saying "It increased my interest to learn using R for analyzing my research data in the future", "It opened my interest to learn more", "encouraged me to learn and practice in the future", or "I am now more interested in really learning R" whereas others alluded to the complexity of learning R saying "I realized a need to learn a lot to get to use R for my purpose" or "I still have a long ways to go but I learned new skills I will apply today."



Figure 3.1. Relative frequency of Likert scale responses for each survey for survey questions related to participant interest in R including questions on intent to use R, usefulness of R, and satisfaction with R. Likert responses are 1 =Strongly disagree, 2 =Disagree, 3 =Somewhat Disagree, 4 =Neutral, 5 =Somewhat Agree, 6 =Agree, and 7 =Strongly Agree.

Figure 3.1 (cont'd)



To evaluate changes in participant confidence learning R, we examine results from Likert statements in the following groupings: ease of using R, difficulty using R, ease of learning R, learning needs, and perceived ability to complete tasks in R. Figure 3.2 summarizes Likert responses in each of these groups, comparing Survey 1 and Survey 2. Results on confidence learning R are mixed, as some groupings show a general increase in confidence such as ease of using R and ability to complete tasks in R, whereas other groupings have uncertain trends between Survey 1 and Survey 2. Some participants expressed an increase in confidence in their qualitative feedback saying "exposure and exercises lower barriers of trying out R" or "the workshop showed me we do not have to know everything by heart and increases my confidence

to use R."



Figure 3.2. Relative frequency of Likert scale responses for each survey for survey questions related to participant confidence in learning R including questions on ease of and difficulty using R, ease of learning R, R learning needs, and ability to complete tasks in R. Note that lower values for difficulty using R correspond with less difficulty using R. Likert responses are 1 = Strongly disagree, 2 = Disagree, 3 = Somewhat Disagree, 4 = Neutral, 5 = Somewhat Agree, 6 = Agree, and 7 = Strongly Agree.

Figure 3.2 (cont'd)



Surveys 2 and 3 also included workshop feedback questions to gather participant opinions on the appropriateness, accessibility, and inclusivity of the workshop. When asked about the most valuable aspects of the seminar, participants frequently mentioned the introduction for beginning R users, error handling and solutions, and group exercises. However, some MSU participants expressed that they would have preferred a slower pace with more detailed explanations, stating, "I think doing less and breaking down the smaller steps more would have been more valuable" and "It is still too advanced for people who are completely new to R." Taking this feedback into account, adjustments were made in the Ethiopia workshop to expand the allotted time for exercises and make them more suitable for beginners. However, participant views on the appropriateness of the material for beginners remained mixed. While some participants found the material too challenging for beginners, offering general feedback such as "It does not work for beginners" or "presenters should understand that all participants are not at the same level with this software," others provided specific recommendations such as:

- "Introduce the basics, such as variable naming and assigning, for those with limited prior experience with R."
- "Make sure that the basic concepts of R, functions, and the like are covered for beginners, considering that participants have different degrees of experience with R."
- "For beginners, start from RStudio with symbols."
- "Strengthen the basics of R for newbies."

On the other hand, some participants found the material appropriate, stating, "It was designed for people with little to no experience with R, so I think it was good for the intended audience/application" and "The workshop is very good for attendees with limited knowledge to get started with data wrangling, etc." In response to this feedback, extended office hours were offered to participants in Ethiopia to address R-related questions, and future iterations of the workshop will expand the "Getting Started with R" exercise script to include more basic concepts such as variable naming, assigning, and commonly used symbols.

Survey feedback also addressed the accessibility and inclusivity of the workshop. Participants appreciated the open-access nature of the materials, stating, "It's accessible to me because I have access to the materials and practical exercises" or "I like that I can review the content at a later time." Regarding inclusivity, participants often mentioned the support provided through the in-person and collaborative aspects of the seminar, with comments such as:

• "The trainer was available to me when I needed assistance."

- "I felt comfortable participating and asking questions."
- "As group members assist beginners, it is inclusive for me."
- "The instructor asked for feedback and was ready to assist."
- "Good workshop atmosphere."

Some participants also appreciated the use of agricultural data, as it aligned with their research area. However, opinions on the appropriateness for beginners varied. While some participants stated, "The workshop is inclusive to me as it starts from the beginning, considering I am a beginner using R" and "I had basic knowledge of R, and this workshop includes beginners like me," others felt that inclusivity could be improved by "matching the needs of participants with different levels of experience in R" and by starting with RStudio, as one participant mentioned, "I am a beginner, so I expected starting with the symbols used."

Based on this feedback, we conclude that the workshop environment was generally accessible and inclusive, fostering a positive and collaborative learning atmosphere. This observation is supported by the quantitative feedback on accessibility and inclusivity presented in Table 3.1, where the majority of participants agreed that the seminar was accessible and inclusive, and would recommend it to others. However, it is important to note that not all participants found every aspect equally accessible and inclusive, leaving room for improvement. For future iterations, one concrete improvement based on participant feedback will be to arrange participants into groups for cooperative exercises, matching beginners with those who have more experience with R, allowing participants to learn from each other more effectively.
Statement	Strongly	Disagree	Somewhat	Neutral	Somewhat	Agree	Strongly
	Disagree		Disagree		Agree		Agree
The seminar portion	0	0	1	0	5	12	9
of this workshop was							
accessible to me.							
The seminar portion	1	0	1	0	3	12	10
of this workshop was							
inclusive to me.							
I would recommend	0	0	0	0	3	9	15
this workshop to							
others.							

Table 3.1. Number of responses for each level of the three feedback Likert statements for the 27 participants completing these statements. Participants generally agreed that the seminar portion of the workshop was accessible and inclusive, and all participants agreed that they would recommend the workshop to others.

3.5 A novel training experience on the theory and application of Bayesian statistics in agronomy for research professionals in Africa

3.5.1 Introduction and workshop motivation

Chapter 2 showed the utility of Bayesian statistical methods for agronomy problems while this chapter gives a first step for introducing researchers to the computational tools necessary to perform Bayesian analysis. As Bayesian methods become more ubiquitous throughout a variety of disciplines, the importance of training opportunities for the theory and application of these methods becomes increasingly important. Furthermore, it is crucial to target specific audiences that would most benefit from the use of Bayesian methods to tackle agronomy problems. Commissioned by the Director of the Sustainable Agrifood Systems program at the International Maize and Wheat Improvement Center (CIMMYT), we offered a novel training experience on the theory and application of Bayesian statistics in agronomy for research professionals in Africa, hosted at the International Livestock Research Institute (ILRI) in Addis Ababa, Ethiopia.

Although CIMMYT headquarters is based in Mexico City, Mexico, the workshop took

place in Ethiopia to target research groups in Africa that could benefit from the use of Bayesian methods in agriculture. This choice was based on the importance of agriculture in Ethiopia and other African countries, as well as the limited opportunity for statistical workshops in African countries. For example, in Ethiopia, agriculture accounts for over a third of the total GDP of the country, and approximately 70% of the Ethiopian workforce works in the agricultural sector (Ayele, 2022). Comparatively, American farms represent approximately 0.7% of the GDP of the United States and only 1.3% of the workforce (Kassel et al., 2023). Beyond the differences in importance of agriculture, the realities of agricultural practices are vastly different in Ethiopia than in the United States and other Western countries. Compared to more developed countries, agriculture in Ethiopia is characterized by small, fragmented plots and a lack of mechanization with 86% of landholding households owning less than 2 hectares of land (Wendimu, 2021) with almost 95% of available farm power coming from human and animal power (Ayele, 2022). These differences in agriculture make it increasingly important for agronomy research on Ethiopian agriculture to take place in Ethiopia, where researchers have a more intimate knowledge of the specific constraints related to agricultural practices in the country. However, in order for researchers in Africa to perform appropriate statistical analyses on agronomy problems, they need access to statistical training opportunities such as workshops. For example, the worldwide organization R-ladies, focused on promoting gender diversity in the R learning community, has 218 global chapters in 29 countries. However, there are only 14 chapters in all of Africa compared to over 50 chapters in the United States alone and over 50 chapters in Europe. This disparity in access to statistical training makes offering training opportunities in African countries even more vital. For these reasons, we began development of a 5-day novel training experience on the fundamentals of Bayesian methods in agronomy to be hosted in

Ethiopia and offer training for researchers throughout Africa.

3.5.2 Workshop overview and agenda

The workshop was designed to introduce early-career agronomy researchers to the fundamentals of the theory and practical application of Bayesian methods and offer a supervised opportunity for researchers to use Bayesian methods on their own data. We had 124 interested applicants and accepted about 3 dozen applicants to maximize impact while also ensuring a small enough group to offer hands-on work based on the number of instructors. The workshop was hosted at ILRI, a research institute co-hosted by the governments of Ethiopia and Kenya that houses more than a dozen international agricultural research and development institutes, making it an ideal location for a broad-impact workshop on Bayesian statistics in agronomy. Participants were selected to maximize the number of agronomy research groups impacted by this work by prioritizing accepting at least one representative participant from as many research centers as possible. We ultimately had 33 participants attend the workshop representing CIMMYT centers in Ethiopia, Zimbabwe, Malawi, and India, the Ethiopian Ministry of Agriculture, the Ethiopian Institute of Agricultural Research, the Ethiopian Agricultural Transformation Institute, the Digital Green Foundation, Hawassa University, the Zimbabwe Ministry of Lands, Agriculture, Fisheries and Rural Development, the Amhara and Gondar Agricultural Research Institutes, and Ethiopian Agricultural Research Centers including Debre Birhan, Debre Markos, Jimma, and Kulumsa.

Local organization of the workshop was led by Dr. Gerald Blasch, Crop Disease Geo-Spatial Data Scientist at CIMMYT, while the workshop agenda and content organization was led by Dr. Frederi Viens, Professor of Statistics at Rice University, and supported by a team of instructors, with the main goal of providing the knowledge and practical skills necessary for

participants to perform their own Bayesian analyses in agronomy research. To facilitate this aim, the workshop was scheduled for five days and structured as follows:

Day 1: Introduction and R Primer. Day 1 began with a 90-minute overview of the workshop and Bayesian statistics in agronomy, presented by Professor Frederi Viens, to motivate the workshop and the use of Bayesian methods for agronomy problems. The remainder of Day 1 included R training to introduce participants to the R programming language, data manipulation and visualization techniques, and error handling, to serve as the computational foundation for the remainder of the week. This was presented by myself and encompassed the seminar portion of the "You Can Learn R" seminar and concluded with the dissemination of the "You Can Learn R" written document, as mentioned in Section 3.4, Presentation and evaluation of results.

Day 2: Theory of Bayesian Statistics. Day 2 was focused on introducing first principles of probability theory applied to Bayesian inference. Early topics included random variables, probability distributions, and basic linear regression, to facilitate leading into more specifics of Bayesian methods including the distinctions between frequentist and Bayesian methods, the advantages of Bayesian methods, and the importance of prior, likelihood, and posterior distributions in the Bayesian framework. This content, though theoretically complex at times, is vital for understanding and interpreting Bayesian analysis in a real-world context. Day 2 instruction was led by Professor Dennis Ikpe, assistant professor in the Department of Statistics and Probability at Michigan State University, and supported by Professor Frederi Viens. **Day 3: Practical application of Bayesian statistics.** Day 3 moved into the computational aspects of using Bayesian methods. Bayesian statistics relies heavily on computation via Markov-chain Monte Carlo (MCMC) for sampling from posterior distributions. To help participants understand the motivation and implementation of Bayesian methods

computationally, Day 3 discussed topics such as conjugate prior distributions, the basic Gibbs sampler, and the practicality of performing more advanced MCMC methods such as Hamiltonian Monte-Carlo using Stan, a platform for performing Bayesian inference, and rstanarm, an interface between R and Stan. This day was designed to provide a fundamental understanding of how Bayesian inference is performed in practice and how to interpret the complex results provided by these computational interfaces, based on the theoretical understanding provided by Day 2. Day 3 instruction was led by Professor Leonard Johnson, teaching specialist in the Department of Statistics and Probability at Michigan State University.

Day 4: Examples of Bayesian statistics in Agronomy. Day 4 provided real-world examples of Bayesian analysis applied to agronomy research. Professor Viens began the discussion with an overview of two of his recent publications that utilize Bayesian analysis to tackle agricultural problems in Malawi. Using these papers as a guide, Viens was able to demonstrate the practical utility of using Bayesian analysis and show how Bayesian results are interpreted in an agronomy context. This high-level overview of applied Bayesian methods gave participants an idea of the types of agronomy questions that can be answered using Bayesian analysis and the final product produced by such analyses. Professor Innocensia John, an economist with the Department of Agricultural Economics and Business at the University of Dar es Salaam, continued the conversation with a presentation of her ongoing work with Viens using Bayesian statistics for agronomy problems in Malawi, including technical coding details related to data-preparation and performance and analysis of Bayesian regression. Finally, I presented my work on agricultural risk mitigation discussed in Chapter 2 as well as a simplified example of Bayesian linear regression with coding details for participants to follow along. We concluded Day 4 with participants forming groups and deciding on their agronomy questions to perform analysis on

their own data in Day 5.

Day 5: Supervised group projects. On Day 5, participants were able to use the knowledge and skills gained in the first four days of the workshop to work collaboratively with other attendees and begin a Bayesian analysis on their own agronomy data. Participants used the coding structure presented in Day 4 as a model to scaffold their Bayesian regression analysis, including implementing informed prior distributions and analyzing results. While working, groups were able to ask questions of the facilitating instructors, myself and Professor John. The workshop concluded with each group presenting the progress made on their analysis to the larger group.

3.5.3 Feedback, impacts, and future work

To evaluate the success of this novel workshop and inform improvements for future iterations, participants were asked to complete a final survey to give feedback on their experience over the five-day workshop. This brief survey included questions about which aspects of each day participants found most valuable and where they saw room for improvement, as well as questions pertaining to their reactions to the workshop overall. Feedback from participants, both verbally throughout the workshop and in the final feedback survey, was largely positive, with some recommendations for future improvement.

In general, participants found the practical exercises and examples most beneficial, such as those offered during the Day 1 R workshop or the examples of Bayesian methods in agronomy and supervised group projects in Day 4 and 5. Constructive feedback from participants discussed how some of the more technical details of Bayesian theory and computation in Day 2 and 3 could have been condensed and supplemented with more hands-on examples. Some participants suggested changes for improvement including extending the supervised group project portion and moving the Day 4 agronomy examples to earlier in the workshop to better motivate the use

of Bayesian methods in agronomy. Overall, participants appreciated this one-of-a-kind training experience and the opportunity to learn practical skills for applying Bayesian methods to their own research. This feedback is best summarized in the following statement from a participant:

"Most important, I liked the workshop very much and I am very grateful for this training opportunity. With [aforementioned] suggestions, I can imagine that the learning experience would be improved considerably. Asking around other participants, most agreed that theory on Day 2 and Day 3 was too much and some participants were scared off and best learning experience were Day 1, Day 4 and Day 5."

This positive feedback coupled with the large number of applicants demonstrates the necessity of such a workshop. Based on positive feedback from participants, organizers, and instructors, we intend to revise and repeat the workshop in future years. Revisions would include more hands-on exercises to aid understanding of the fundamentals of Bayesian methods as well as additional time and focus for real-world examples of using Bayesian statistics in agronomy and supervised group projects. Further, we would like to provide some workshop materials in advance to maximize the effectiveness of the in-person workshop and perform some follow-up with participants to offer continued support they use Bayesian methods in their own work. Finally, we hope to expand this effort in the long-term to offer similar workshops in other African countries and to train local researchers in Africa to be able to offer such training experiences in the future. Ultimately, this novel training experience has served as a starting point for making the use of Bayesian statistics in agronomy more widespread, particularly for countries where agricultural research is so vital, like Ethiopia.

3.6 Limitations and future work

The "You Can Learn R" workshop was renewed for the Summer 2023 College of Natural

Science Great IDEA Fellowship at MSU with the intention of expanding the content and impact of the workshop. While I was able to create an initial iteration of the workshop for presentation at MSU and in Ethiopia, my plan is to continue refining the workshop materials to improve accessibility and inclusivity and to broaden its reach.

The initial iterations were limited to in-person workshop environments with under 50 total participants. After refining the workshop, my goal is to present it again to the MSU community, incorporating a hybrid option to reach a wider audience or adjusting the materials to an asynchronous format. This would allow interested individuals to engage with the workshop materials at their own pace. When advertising my workshop at MSU, I quickly reached the maximum capacity of 32 participants within 24 hours of announcing it, with a total of 99 participants signing up. The demand for such a workshop is evident, and I aim to expand and refine my work to reach a larger audience.

Evaluation of the workshop's effectiveness in terms of participant interest, confidence, and ability in learning R, as well as the accessibility and inclusivity of the resource, was hindered by both the limited number and quality of survey responses. For instance, Survey 3 only had 5 participants from the iteration in Ethiopia, and some feedback responses were confounded with those related to the larger Bayesian statistics workshop conducted in Ethiopia, rather than solely focusing on the "You Can Learn R" workshop. These limitations made it challenging to accurately assess the effectiveness of the written resource and working sessions due to the reliance on a limited set of survey results.

Regarding the measurement of interest, confidence, and ability in learning R, the survey questions, which were adapted from established user experience surveys, did not consistently align with the three desired dimensions. While Likert statements generally addressed interest and

confidence in learning R, assessing ability proved to be more complex. The only measure of ability relied on qualitative questions regarding participants' completion of tasks, making it impractical to track changes in ability over time. In future iterations, I intend to refine the survey questions to more accurately capture the three dimensions of interest, confidence, and ability in learning R. Additionally, I will incorporate the collection of completed exercises to gain a clearer understanding of the workshop's impact on participants' ability to learn R.

Feedback on accessibility and inclusivity was also limited due to the design of the survey questions. To address this limitation, my plan is to revise the questions related to accessibility and inclusivity by separating them to solicit distinct responses regarding both positive and negative aspects. For example, instead of combining the questions as, "In what ways was this workshop accessible to you? In what ways was it not?", I will ask them separately to ensure participants provide comprehensive answers to each aspect. Moreover, I will take steps to clarify the meaning of the terms "accessible" and "inclusive", as some participants expressed uncertainty when asked about the inclusivity of the workshop. Furthermore, I intend to enhance the effectiveness of the feedback survey questions and supplement the data with one-on-one interviews with selected future participants. This approach will provide a more complete qualitative view of participants' perceptions of the workshop. By improving the workshop and its feedback mechanisms, "You Can Learn R" will continue to evolve as a dynamic workshop, connecting R learners both at MSU and around the world.

3.7 Conclusion

The R programming language has experienced significant growth in terms of the number of R packages available and the diverse range of fields utilizing the language. Its flexibility, versatility, and open-source nature have made it an invaluable tool for researchers at all levels,

enabling them to perform essential data-driven tasks such as data manipulation, visualization, and communication. However, the extensive array of R packages and learning resources, coupled with the steep learning curve, can overwhelm researchers and hinder their adoption of R for their own work. This is particularly challenging for early-career researchers, who face time constraints and may find it difficult to commit to lengthy learning materials like textbooks or courses lasting multiple months.

To address this need for a resource that facilitates researchers in quickly grasping R and leveraging its benefits in their work, I developed the "You Can Learn R" workshop. This workshop offers an accessible and inclusive learning experience tailored specifically for earlycareer researchers. It comprises a seminar portion to introduce participants to R and foster cooperative learning, a comprehensive written document to accompany their R learning journey, and a hands-on working session in a collaborative environment, enabling participants to apply R to their own projects without fear of making mistakes.

The initial pilot iterations of the workshop were conducted with MSU graduate students and postdoctoral researchers, as well as part of a larger workshop on Bayesian methods for earlycareer agronomy professionals in Africa. These pilots effectively demonstrated the necessity of such a resource and the workshop's efficacy. Participants in both settings found the workshop to be accessible and inclusive, appreciating such aspects as the provision of open-access materials and the encouragement of peer-based cooperative learning through group exercises. Nevertheless, valuable feedback received highlighted areas for improvement, such as pairing participants with varying levels of R experience and covering more foundational R concepts for beginners.

To further refine and expand the workshop's impact, the project has received support

from the MSU College of Natural Science's Great IDEA fellowship. This ongoing support will enable continuous improvement based on participant feedback and evolving trends in R learning. The workshop will continue to serve as a valuable starting point and companion resource, aiding researchers from diverse backgrounds in acquiring R skills and reaping the benefits of this powerful programming language.

CHAPTER 4: CONCLUSION

Bayesian methods are gaining increasing popularity across disciplines due to their numerous advantages. Throughout my graduate career at MSU, I have acquired the fundamental knowledge in theory and computation necessary for performing Bayesian statistics, as well as the teaching skills required to make Bayesian methods accessible to researchers who can benefit from their advantages. Through three research projects focused on the practical realities of performing and teaching Bayesian methods, I have come to understand that being a Bayesian statistician also entails being a statistics educator.

As a statistics educator, I believe that statistics is relevant to varying degrees for everyone, and the teaching of statistics should be tailored to the specific needs of the audience. This need is particularly apparent in Bayesian statistics, as using Bayesian methods or even understanding the results of Bayesian analysis requires a certain level of understanding of probability theory and computational ability and resources.

One key advantage of Bayesian statistics is its intuitive probabilistic interpretations, robust uncertainty quantification, and the ability to incorporate domain-specific information through prior distributions. However, realizing these theoretical advantages necessitates a solid foundational understanding of probability and Bayesian methods in order to effectively use and comprehend them. Therefore, it is the responsibility of Bayesian statisticians to possess their own understanding of the theoretical underpinnings of Bayesian statistics and to effectively communicate these foundations to other researchers and stakeholders.

In Chapter 2, we employed Bayesian methods to quantify the risk mitigation effect associated with increased rotational complexity in farming systems in the Midwest US. This project required a strong theoretical understanding of Bayesian modeling to conduct Bayesian

regression and effectively communicate the results. Prior to conducting the Bayesian analysis, we developed a modeling plan and presented it to both Land Core, the non-profit organization leading the project, and Compeer Financial, a Midwest farm lender supporting the work. Throughout the modeling process, we provided explanations of the methodology and potential outcomes to the entire multidisciplinary team, and we generated regular reports to update Compeer Financial. Once we generated preliminary model results, I presented our work as an oral presentation at the Conference on Applied Statistics in Agriculture and Natural Resources, sharing our progress and potential impact with other researchers in the field. Additionally, when applying for grants to fund our work, it was essential to concisely and convincingly convey our modeling process and preliminary results to grant reviewers. Furthermore, during the development of the upcoming tool for farm lenders and insurers, it was crucial to effectively communicate the results to user interface developers so that they could accurately represent the findings visually. As our project team continues to expand, teaching the theoretical foundations of our work remains a priority to bring new members up to speed.

Throughout each step of the project, it was vital to adapt my communication approach to suit the knowledge and needs of different audiences. For instance, the chief strategy officer at Land Core, who had limited statistical knowledge, was responsible for conveying our analysis results to audiences such as potential funders and government policy-makers. It was imperative that I communicated the necessary theoretical understanding to the chief strategy officer so that they could accurately convey the results and impact of our work, thereby garnering support for our research and influencing federal policies related to soil practices.

Conducting Bayesian analysis requires not only an understanding of probability theory but also the ability to effectively communicate the theoretical knowledge required for each

audience impacted by the project. Consequently, the communication of the theoretical basis for Bayesian methods was a critical aspect of teaching agronomists in Africa how to perform Bayesian analysis, as discussed in Section 3.5. To fully leverage the benefits of Bayesian analysis, including the use of informed priors and the interpretation of probabilistic results, an understanding of the relationship between prior, likelihood, and posterior distributions is necessary. To grasp this relationship, a foundational knowledge of probability theory, including concepts such as random variables, probability distributions, and Bayes' theorem, is essential. Through our innovative workshop on teaching the fundamentals of Bayesian statistics in agronomy, we aimed to share the necessary theoretical knowledge for participants to conduct their own applied Bayesian analysis in agronomy. Although the theory may be challenging at times, we facilitated comprehension through practical exercises and step-by-step examples, equipping participants with the theoretical knowledge required to perform Bayesian analysis effectively.

Beyond theoretical concerns, Bayesian methods also require the practitioner to have computational ability and resources. Since Bayesian posterior distributions often cannot be computed explicitly, posterior distributions must be generated using Markov Chain Monte Carlo (MCMC) methods. These methods are computationally intensive and can generate a large amount of output data if storing full or even partial posterior distributions is a priority. In our analysis discussed in Chapter 2, being the first of its kind to use Bayesian methods on a field-level agricultural dataset of such magnitude, we encountered significant computational challenges. The initial aggregation of the dataset by our data manager took multiple months and required oversight by all other computationally savvy team members to review the aggregating code and verify data accuracy. Due to the size and regional variation of our data, substantial data

manipulation was necessary to explore modeling trends on spatial subsets of the data. For Bayesian analysis, we investigated various computational interfaces for using Bayes in R, such as rstan, brms, rstanarm, and NIMBLE, ultimately selecting brms for its flexibility in prior definitions and ability to fit large data in reasonable time. Fitting our Bayesian models required a high-performance computing cluster, and model fitting and prediction took approximately 90 minutes for each county. Storing the resultant prediction data created large CSV files for each county, and managing and analyzing the data from all county-level files together became challenging. As both our input and output data continue to expand, we are transitioning to a comprehensive, cloud-hosted, relational database to provide efficient storage and querying of data as the project progresses.

It is evident that the use of Bayesian methods requires a certain degree of computational understanding. Therefore, a necessary step in learning Bayesian methods is the ability to use a statistical programming language, such as R. In Chapter 3, I presented design considerations for making learning R accessible and inclusive for a variety of audiences. Through the "You Can Learn R" workshop, I provided a starting point for researchers with little to no R experience to learn R and use it for their research. This starting point served as a foundational module for the larger workshop discussed in Section 3.5, providing participants with the computational skills necessary to facilitate Bayesian analysis. Through this novel training experience, participants had collaborative and supervised opportunities to engage with the computational realities of Bayesian statistics.

As Bayesian statistics continues to gain popularity, Bayesian statisticians have a responsibility as statistics educators to make their applied Bayesian research accessible to a wide audience and to facilitate the informed use of Bayesian statistics in other disciplines. Each of the

three projects discussed has demonstrated the complexity and considerations necessary for both performing and teaching Bayesian methods. However, these projects are just the starting point for my continued work as a Bayesian statistician and statistics educator, particularly as I embark on my research associate position at the Center for Statistical Training and Consulting (CSTAT) at MSU.

In Chapter 2, I discussed how we used Bayesian methods to quantify the risk mitigation in terms of corn yield associated with higher rotational complexity, focusing on two states in the Midwestern US. This project is the first step in establishing an empirical link between regenerative soil practices and crop yield, and our work will continue to expand as we include more management practices, crops, weather conditions, and geographical areas. I plan to continue collaborating on this work in a reduced capacity as I begin my upcoming research position.

In Chapter 3, I presented a workshop to teach researchers how to learn R and the fundamentals of data analysis. This workshop serves as a foundation for researchers to develop their computational skills and apply them to various statistical methodologies, including Bayesian analysis. I envision expanding this workshop to reach a broader audience and continuing to add resources to make it accessible to a wider range of individuals interested in learning R for their own research. Further, the "You Can Learn R" workshop served as a foundation for the workshop we developed to teach agronomists in Africa the fundamentals of Bayesian statistics, detailed in Section 3.5. This workshop aimed to empower agronomists with the knowledge and skills to apply Bayesian methods to their research and decision-making processes. Moving forward, I plan to continue collaborating on future iterations of this workshop, bridging the gap between statistical theory and its practical application to agronomy,

particularly for scholars in African countries.

As I take on my research associate position at CSTAT, my primary focus will be on providing statistical training and consulting services to researchers across various disciplines. This role offers an opportunity to further contribute to the advancement and application of Bayesian methods by working closely with researchers and helping them integrate Bayesian analysis into their research projects. Additionally, I plan to actively engage in outreach efforts, organizing workshops and seminars to promote the understanding and adoption of Bayesian statistics among the research community.

In conclusion, being a Bayesian statistician entails not only possessing a strong theoretical understanding of Bayesian methods but also being an effective statistics educator. The ability to communicate complex statistical concepts to different audiences, adapt teaching strategies to diverse learners, and provide computational guidance are crucial aspects of successfully applying and teaching Bayesian statistics. By combining my expertise in Bayesian methods, computational skills, and statistics education, I am committed to advancing the field of Bayesian statistics and empowering researchers to utilize these powerful methods in their own work.

BIBLIOGRAPHY

- Abatzoglou, J. T., Dobrowski, S. Z., Parks, S. A., & Hegewisch, K. C. (2018). TerraClimate, a high-resolution global dataset of monthly climate and climatic water balance from 1958– 2015. *Scientific Data*, 5(1), 170191. https://doi.org/10.1038/sdata.2017.191
- Albers, M. A., Dobos, R. R., & Robotham, M. P. (2022). User Guide for the National Commodity Crop Productivity Index (NCCPI) Version 3.0. USDA National Resources Conservation Service, Soil and Plant Science Division.
- Allaire, J. J., Xie, Y., Dervieux, C., McPherson, J., Luraschi, J., Ushey, K., Atkins, A., Wickham, H., Cheng, J., Chang, W., & Iannone, R. (2023). *rmarkdown: Dynamic Documents for R*. https://github.com/rstudio/rmarkdown
- Ayele, S. (2022). The resurgence of agricultural mechanisation in Ethiopia: Rhetoric or real commitment? *The Journal of Peasant Studies*, 49(1), 137–157. https://doi.org/10.1080/03066150.2020.1847091
- Beal Cohen, A. A., Seifert, C. A., Azzari, G., & Lobell, D. B. (2019). Rotation Effects on Corn and Soybean Yield Inferred from Satellite and Field-level Data. *Agronomy Journal*, 111(6), 2940–2948. https://doi.org/10.2134/agronj2019.03.0157
- Boryan, C., Yang, Z., Mueller, R., & Craig, M. (2011). Monitoring US agriculture: The US Department of Agriculture, National Agricultural Statistics Service, Cropland Data Layer Program. *Geocarto International*, 26(5), 341–358. https://doi.org/10.1080/10106049.2011.562309
- Bowles, T. M., Mooshammer, M., Socolar, Y., Calderón, F., Cavigelli, M. A., Culman, S. W., Deen, W., Drury, C. F., Garcia Y Garcia, A., Gaudin, A. C. M., Harkcom, W. S., Lehman, R. M., Osborne, S. L., Robertson, G. P., Salerno, J., Schmer, M. R., Strock, J., & Grandy, A. S. (2020). Long-Term Evidence Shows that Crop-Rotation Diversification Increases Agricultural Resilience to Adverse Growing Conditions in North America. *One Earth*, 2(3), 284–293. https://doi.org/10.1016/j.oneear.2020.02.007
- Brooke, J. (1995). SUS: A quick and dirty usability scale. Usability Eval. Ind., 189.
- Bürkner, P.-C. (2017). **brms**: An *R* Package for Bayesian Multilevel Models Using *Stan. Journal* of *Statistical Software*, 80(1). https://doi.org/10.18637/jss.v080.i01
- Cassman, K. G., & Grassini, P. (2020). A global perspective on sustainable intensification research. *Nature Sustainability*, *3*(4), 262–268. https://doi.org/10.1038/s41893-020-0507-8
- Deines, J. M., Guan, K., Lopez, B., Zhou, Q., White, C. S., Wang, S., & Lobell, D. B. (2023). Recent cover crop adoption is associated with small maize and soybean yield losses in the United States. *Global Change Biology*, 29(3), 794–807. https://doi.org/10.1111/gcb.16489

- Deines, J. M., Patel, R., Liang, S.-Z., Dado, W., & Lobell, D. B. (2021). A million kernels of truth: Insights into scalable satellite maize yield mapping and yield gap analysis from an extensive ground dataset in the US Corn Belt. *Remote Sensing of Environment*, 253, 112174. https://doi.org/10.1016/j.rse.2020.112174
- Deines, J. M., Wang, S., & Lobell, D. B. (2019). Satellites reveal a small positive yield effect from conservation tillage across the US Corn Belt. *Environmental Research Letters*, *14*(12), 124038. https://doi.org/10.1088/1748-9326/ab503b
- Dogucu, M., Johnson, A. A., & Ott, M. (2023). Framework for Accessible and Inclusive Teaching Materials for Statistics and Data Science Courses. *Journal of Statistics and Data Science Education*, 1–7. https://doi.org/10.1080/26939169.2023.2165988
- Dunson, D. B. (2001). Commentary: Practical Advantages of Bayesian Analysis of Epidemiologic Data. *American Journal of Epidemiology*, *153*(12), 1222–1226. https://doi.org/10.1093/aje/153.12.1222
- Feng, Z., Leung, L. R., Hagos, S., Houze, R. A., Burleyson, C. D., & Balaguru, K. (2016). More frequent intense and long-lived storms dominate the springtime trend in central US rainfall. *Nature Communications*, 7(1), 13429. https://doi.org/10.1038/ncomms13429
- Gallagher, J. (2022, October 5). Learn R: Best Courses, Books, and Resources for Learning R. *Career Karma*. https://careerkarma.com/blog/how-to-learn-r/
- Jin, Z., Azzari, G., & Lobell, D. B. (2017). Improving the accuracy of satellite-based highresolution yield estimation: A test of multiple scalable approaches. *Agricultural and Forest Meteorology*, 247, 207–220. https://doi.org/10.1016/j.agrformet.2017.08.001
- Kang, Y., & Özdoğan, M. (2019). Field-level crop yield mapping with Landsat using a hierarchical data assimilation approach. *Remote Sensing of Environment*, 228, 144–163. https://doi.org/10.1016/j.rse.2019.04.005
- Kassel, K., Lanigan, T., Martin, A., Michael-Midkiff, J., Russell, D., Ruth, T., Sanguinett, C., Smits, J., Symanski, E., Kassel, K., Lanigan, T., Martin, A., Michael-Midkiff, J., Russell, D., Ruth, T., Sanguinett, C., Smits, J., & Symanski, E. (2023). Selected Charts from Ag and Food Statistics: Charting the Essentials, February 2023. https://doi.org/10.22004/AG.ECON.333548
- Kimm, H., Guan, K., Gentine, P., Wu, J., Bernacchi, C. J., Sulman, B. N., Griffis, T. J., & Lin, C. (2020). Redefining droughts for the U.S. Corn Belt: The dominant role of atmospheric vapor pressure deficit over soil moisture in regulating stomatal behavior of Maize and Soybean. Agricultural and Forest Meteorology, 287, 107930. https://doi.org/10.1016/j.agrformet.2020.107930
- Kravchenko, A. N., Snapp, S. S., & Robertson, G. P. (2017). Field-scale experiments reveal persistent yield gaps in low-input and organic cropping systems. *Proceedings of the National Academy of Sciences*, 114(5), 926–931.

https://doi.org/10.1073/pnas.1612311114

- Li, X., Tack, J. B., Coble, K. H., Barnett, B. J., Li, X., Tack, J. B., Coble, K. H., & Barnett, B. J. (2016). Can Crop Productivity Indices Improve Crop Insurance Rates? https://doi.org/10.22004/AG.ECON.235750
- Lobell, D. B., Roberts, M. J., Schlenker, W., Braun, N., Little, B. B., Rejesus, R. M., & Hammer, G. L. (2014). Greater Sensitivity to Drought Accompanies Maize Yield Increase in the U.S. Midwest. *Science*, 344(6183), 516–519. https://doi.org/10.1126/science.1251423
- Lobell, D. B., Thau, D., Seifert, C., Engle, E., & Little, B. (2015). A scalable satellite-based crop yield mapper. *Remote Sensing of Environment*, 164, 324–333. https://doi.org/10.1016/j.rse.2015.04.021
- Lund, A. (2001). Measuring Usability with the USE Questionnaire. Usability and User Experience Newsletter of the STC Usability SIG, 8.
- Marini, L., St-Martin, A., Vico, G., Baldoni, G., Berti, A., Blecharczyk, A., Malecka-Jankowiak, I., Morari, F., Sawinska, Z., & Bommarco, R. (2020). Crop rotations sustain cereal yields under a changing climate. *Environmental Research Letters*, 15(12), 124011. https://doi.org/10.1088/1748-9326/abc651
- Natural Resources Conservation Service, U. S. D. O. A. (2016). Gridded Soil Survey Geographic Database (gSSURGO) [dataset]. Natural Resources Conservation Service, United States Department of Agriculture. https://doi.org/10.15482/USDA.ADC/1255234
- Newton, J. (2019, October 30). *Farm bankruptcies rise again*. Wisconsin State Farmer. https://www.wisfarmer.com/story/news/2019/10/30/farm-bankruptcies-filings-up-24over-year-ago/4096381002/
- Ortiz-Bobea, A., Knippenberg, E., & Chambers, R. G. (2018). Growing climatic sensitivity of U.S. agriculture linked to technological change and regional specialization. *Science Advances*, *4*(12), eaat4343. https://doi.org/10.1126/sciadv.aat4343
- PRISM Climate Group. (2022). PRISM Normals [dataset]. https://prism.oregonstate.edu
- Prost, L., Makowski, D., & Jeuffroy, M.-H. (2008). Comparison of stepwise selection and Bayesian model averaging for yield gap analysis. *Ecological Modelling*, 219(1–2), 66– 76. https://doi.org/10.1016/j.ecolmodel.2008.07.026
- R Core Team. (2021). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. https://www.R-project.org/
- Rodell, M., Houser, P. R., Jambor, U., Gottschalck, J., Mitchell, K., Meng, C.-J., Arsenault, K., Cosgrove, B., Radakovich, J., Bosilovich, M., Entin, J. K., Walker, J. P., Lohmann, D., & Toll, D. (2004). The Global Land Data Assimilation System. *Bulletin of the American Meteorological Society*, 85(3), 381–394. https://doi.org/10.1175/BAMS-85-3-381

- Sanford, G. R., Jackson, R. D., Booth, E. G., Hedtcke, J. L., & Picasso, V. (2021). Perenniality and diversity drive output stability and resilience in a 26-year cropping systems experiment. *Field Crops Research*, 263, 108071. https://doi.org/10.1016/j.fcr.2021.108071
- Seifert, C. A., Azzari, G., & Lobell, D. B. (2018). Satellite detection of cover crops and their effects on crop yield in the Midwestern United States. *Environmental Research Letters*, 13(6), 064033. https://doi.org/10.1088/1748-9326/aac4c8
- Seifert, C. A., Roberts, M. J., & Lobell, D. B. (2017). Continuous Corn and Soybean Yield Penalties across Hundreds of Thousands of Fields. *Agronomy Journal*, 109(2), 541–548. https://doi.org/10.2134/agronj2016.03.0134
- Sheng, Y. P., Paramygin, V. A., Rivera-Nieves, A. A., Zou, R., Fernald, S., Hall, T., & Jacob, K. (2022). Coastal marshes provide valuable protection for coastal communities from storminduced wave, flood, and structural loss in a changing climate. *Scientific Reports*, 12(1), 3051. https://doi.org/10.1038/s41598-022-06850-z
- Socolar, Y., Goldstein, B. R., De Valpine, P., & Bowles, T. M. (2021). Biophysical and policy factors predict simplified crop rotations in the US Midwest. *Environmental Research Letters*, *16*(5), 054045. https://doi.org/10.1088/1748-9326/abf9ca
- Swain, S., & Hayhoe, K. (2015). CMIP5 projected changes in spring and summer drought and wet conditions over North America. *Climate Dynamics*, 44(9–10), 2737–2750. https://doi.org/10.1007/s00382-014-2255-9
- Theobold, A., & Hancock, S. (2019). HOW ENVIRONMENTAL SCIENCE GRADUATE STUDENTS ACQUIRE STATISTICAL COMPUTING SKILLS. *STATISTICS EDUCATION RESEARCH JOURNAL*, *18*(2), 68–85. https://doi.org/10.52041/serj.v18i2.141
- USDA/NASS. (n.d.). USDA/NASS QuickStats Ad-hoc Query Tool. Retrieved July 9, 2023, from https://quickstats.nass.usda.gov/results/A64B0F5C-B26F-3A05-8E58-BCA33B34B566
- USGCRP. (2018). *Fourth National Climate Assessment* (pp. 1–470). U.S. Global Change Research Program, Washington, DC. https://nca2018.globalchange.gov
- Van De Pol, M., & Wright, J. (2009). A simple method for distinguishing within- versus between-subject effects using mixed models. *Animal Behaviour*, 77(3), 753–758. https://doi.org/10.1016/j.anbehav.2008.11.006
- Venkatesh, Morris, Davis, & Davis. (2003). User Acceptance of Information Technology: Toward a Unified View. *MIS Quarterly*, 27(3), 425. https://doi.org/10.2307/30036540
- Wickham, H. (2016). ggplot2: Elegant Graphics for Data Analysis (2nd ed. 2016). Springer International Publishing : Imprint: Springer. https://doi.org/10.1007/978-3-319-24277-4

- Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L., François, R., Grolemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T., Miller, E., Bache, S., Müller, K., Ooms, J., Robinson, D., Seidel, D., Spinu, V., ... Yutani, H. (2019). Welcome to the Tidyverse. *Journal of Open Source Software*, 4(43), 1686. https://doi.org/10.21105/joss.01686
- Wickham, H., François, R., Henry, L., & Müller, K. (2020). A Grammar of Data Manipulation [*R package dplyr version 1.0.2*].
- Wolff, S., Schulp, C. J. E., & Verburg, P. H. (2015). Mapping ecosystem services demand: A review of current research and future perspectives. *Ecological Indicators*, 55, 159–171. https://doi.org/10.1016/j.ecolind.2015.03.016
- Worsley, S. (2022). *What is R? The Statistical Computing Powerhouse*. https://www.datacamp.com/blog/all-about-r
- Xie, Y. (2020). *bookdown: Authoring Books and Technical Documents with R Markdown*. https://github.com/rstudio/bookdown
- Xie, Y., Allaire, J. J., & Grolemund, G. (2018). *R Markdown: The Definitive Guide*. Chapman and Hall/CRC. https://bookdown.org/yihui/rmarkdown
- Xie, Y., Dervieux, C., & Riederer, E. (2020). *R Markdown Cookbook*. Chapman and Hall/CRC. https://bookdown.org/yihui/rmarkdown-cookbook
- Xu, T., Guan, K., Peng, B., Wei, S., & Zhao, L. (2021). Machine Learning-Based Modeling of Spatio-Temporally Varying Responses of Rainfed Corn Yield to Climate, Soil, and Management in the U.S. Corn Belt. *Frontiers in Artificial Intelligence*, 4, 647999. https://doi.org/10.3389/frai.2021.647999
- Yan, L., & Roy, D. P. (2016). Conterminous United States crop field size quantification from multi-temporal Landsat data. *Remote Sensing of Environment*, 172, 67–86. https://doi.org/10.1016/j.rse.2015.10.034
- Yigezu Wendimu, G. (2021). The challenges and prospects of Ethiopian agriculture. *Cogent Food & Agriculture*, 7(1), 1923619. https://doi.org/10.1080/23311932.2021.1923619
- Yu, J., Smith, A., & Sumner, D. A. (2018). Effects of Crop Insurance Premium Subsidies on Crop Acreage. American Journal of Agricultural Economics, 100(1), 91–114. https://doi.org/10.1093/ajae/aax058