

# A Markov chain model of crop conditions and intrayear crop yield forecasting

J. R. Stokes 

Department of Agricultural Economics,  
University of Nebraska-Lincoln, Lincoln,  
Nebraska, USA

## Correspondence

J. R. Stokes, Department of Agricultural  
Economics, University of Nebraska-  
Lincoln, 308B Filley Hall, Lincoln, NE  
68583, USA.

Email: [jeffrey.stokes@unl.edu](mailto:jeffrey.stokes@unl.edu)

## Abstract

Crop condition reports are an important source of information for producers, grain traders, businesses, and policymakers to assess and manage the price and yield risk inherent in a given crop. A Markov chain model is proposed for describing the weekly dynamic behavior of reported crop conditions. Empirical transition probabilities are estimated for corn grown in Nebraska, and forecasted crop conditions from the Markov chain are used as inputs to forecast final crop yields prior to harvest time. The results suggest that the modeling and forecasting approach has value for estimating crop yields as intrayear information about crop conditions materializes.

## KEYWORDS

crop condition, Markov chain, maximum entropy, weighted least squares

## 1 | INTRODUCTION

Crop progress reports summarize the results of weekly crop progress and condition surveys conducted by the National Agricultural Statistics Service, a division of the US Department of Agriculture (USDA). The information from the reports is used by producers, grain traders, and businesses, as well as federal and state agencies, to assess and manage the risk inherent in crop production. Research by Lehecka (2014) shows that the reports have substantial informational value in that market prices tend to react rapidly to new crop condition information. As the name implies, the reports provide information about the growing season progress and overall condition of major US crops. The crop condition portion of each report shows the percent of a given crop rated “very poor,” “poor,” “fair,” “good,” and “excellent” for selected states and the United States

as a whole. According to the USDA–National Agricultural Statistics Service:

All states participate in the survey. Each state maintains a list of reporters, largely extension agents and Farm Service Agency staff, who report progress and conditions of selected crops in their area for the current week. Nearly every county in every state has at least one reporter. Reports returned each week account for over 75 percent of the acreage for major commodities.

While subjective in nature, reported crop conditions do provide useful information about intraseason crop quality. For example, a weighted average crop condition index, constructed using an index of 1 (*very poor*) through 5 (*excellent*) with the crop condition percentages as the weights, shows how the quality of a crop varies week to week during the growing season. When the index

Any errors or omissions are the sole responsibility of the author.

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2023 The Author. *Journal of Forecasting* published by John Wiley & Sons Ltd.

increases (decreases), it is the consensus opinion of survey respondents that a better (worse) crop will ultimately result at the culmination of the growing season. It should be noted that while crop conditions are inherently qualitative, an “excellent” condition, for example, implies above average production. In fact, as noted by Norwood and Fackler (1999), the USDA has specific definitions in mind for each qualitative condition:

- Excellent: Yield prospects are above normal. Crops are experiencing little or no stress. Disease, insect damage, and weed pressures are insignificant.
- Good: Yield prospects are normal. Moisture levels are adequate, and disease, insect damage, and weed pressures are normal.
- Fair: less than normal crop condition. Yield loss is a possibility, but the extent is unknown.
- Poor: heavy degree of loss of yield potential which can be caused by excess soil moisture, drought, disease, and so forth.
- Very poor: extreme degree of loss of yield potential, complete or near crop failure.

Therefore, while the class definitions are purely qualitative in nature, they are meant to convey a sense of the quantity of a crop likely to be produced.

Even so, there is no research connecting reported crop conditions, which occur weekly during the growing season, with the ongoing and eventual total production of a crop. Kruse and Smith (1994) developed a model to estimate corn and soybean yields by crop condition class. Heteroskedasticity is noted as a problem attributable to their use of pooled data (i.e., nonconstant yield variance across states), and the authors use weighted least squares to remediate the problem. Norwood and Fackler (1999), building on the research by Kruse and Smith (1994), use ordinary and weighted least squares to estimate the ratio of crop (corn, cotton, soybean, and spring wheat) yield-to-yield trend by crop condition class citing the same source of heteroskedasticity.

The purpose of this research is to report on an alternative yield forecast model based on crop conditions. More specifically, we specify a model of the intrayear movement of reported crop conditions as a Markov chain. Markov chains have seen extensive use in forecasting by, for example, Liu et al. (2015) for mortgage stress testing, Lo et al. (2016) for latent volatility, Tang et al. (2018) for scenario analysis, and Li and Andersson (2020) for density forecasting. Markov chains specific to forecasting crop yields have been employed by Matis et al. (1985, 1989), but those efforts predate the USDA crop condition data described above. Also, our models differ considerably from previous research in that no attempt is made to forecast intrayear or end-of-year yield by condition class directly. Rather, we

focus on the modeling of weekly dynamics of crop conditions and what those dynamics mean for the ongoing and final estimate of overall crop yield. An empirical application of the model is presented for corn grown in Nebraska using two different estimators for the transition probabilities making up the Markov chain.

The paper is organized as follows: In Section 2, we motivate the Markov chain model for crop conditions and present two models for estimating the transition probabilities of the Markov chain; while in Section 3, we provide empirical support for the model. Section 4 concludes.

## 2 | MODELS OF CROP CONDITION

In this section, we briefly motivate reported crop condition as a stochastic process amenable to estimation as a first-order Markov chain. Further, we present two methods of estimation and show how the resulting models can be used to estimate year-end crop yields.

### 2.1 | Crop conditions

Let  $a_i(t)$  denote the number of acres of a given crop (e.g., corn) in a county in the  $i$ th crop condition class at time  $t$  where  $i = 1, \dots, 5$  and  $\sum_i a_i(t) = A(t)$  where  $A(t)$  is the total acreage of a particular crop in the county. Although it is of no consequence for the discussion that follows, to be consistent with constructed crop condition indices, we assume  $i = 1$  represents the “very poor” crop condition,  $i = 2$  represents the “poor” crop condition, and so forth.

The acreage,  $a_i(t)$ , indicates the sequence of crop condition states frequented at  $t = 0, \dots, T$  points in time. If the  $a_i(t)$  were directly observable, it would be straightforward to calculate  $a_{ij} = \sum_t a_{ij}(t)$  where  $a_{ij}(t)$  is the number of acres transitioning from crop condition class  $i$  to  $j$  over one (weekly) time period. Then, assuming a time homogeneous Markov chain, the distribution of the  $a_{ij}(t)$  may be obtained by considering the  $a_i(t-1) = \sum_j a_{ij}(t)$  observations on a multinomial distribution with transition probabilities,  $p_{ij}$ . Anderson and Goodman (1957) have shown that the maximum likelihood estimator of the transition probabilities is then  $p_{ij} = a_{ij} / \sum_j a_{ij} \geq 0$ .

Unfortunately, the  $a_i(t)$  are not directly observable (at least publicly), and there is no guarantee that the underlying Markov chain is time homogeneous. For example, it is reasonable that transition probabilities between crop condition classes may depend on those elements used to define and describe the classes, namely, the amount of soil moisture and the extent of disease and insect damage. As a result, we take the following approach. Let  $c_i(t)$

represent the reported percentage of a crop in the  $i$ th crop condition class at time  $t$  in a given area where for any time  $t$ ,  $\sum_i c_i(t) = 1$ . It follows that observable  $c_i(t) = a_i(t)/A(t)$  are generated from the unobservable  $a_i(t)$ .<sup>1</sup>

When all we can observe are the proportions  $c_i(t)$ , Lee et al. (1970) has shown that there are no set of transition probabilities that will satisfy the dynamic relationship between the unconditional probabilities of the Markov chain.<sup>2</sup> However, if we substitute the crop condition proportions for the unconditional probabilities and allow for random error, denoted  $\tilde{\epsilon}_j(t)$ , the Markov relation may be stated as

$$c_j(t) = \sum_i p_{ij}(t) c_i(t-1) + \tilde{\epsilon}_j(t), \quad \forall j, t, \quad (1)$$

from which the  $p_{ij}(t)$  may be estimated from the  $c_i(t)$  and  $c_j(t-1)$  data.

Equation (1) is the central assumption behind the present research, namely, that a first-order Markov chain (with error) is a reasonable model of crop condition dynamics for a given crop in a given region. Further, notice that in (1), the transition probabilities are shown as a function of  $t$  indicating a time inhomogeneous Markov chain is not precluded. In the empirical section of the paper below, we conduct extensive statistical tests to determine whether there is empirical support for a time inhomogeneous or time homogeneous Markov chain (i.e., whether  $p_{ij}(t) = p_{ij}$ ).

## 2.2 | Transition probability models

Two models are presented in this section for the estimation of the transition probabilities. Tests of the stationarity of a Markov chain have been proposed by Anderson and Goodman (1957), Kelton and Kelton (1984), and Mattsson and Thorburn (1989). As shown in the empirical section below, using the approaches suggested by Kelton and Kelton (1984), there is no statistical evidence in support of a time inhomogeneous Markov chain for corn crop conditions in Nebraska. Therefore, we proceed with the model presentations and descriptions assuming that  $p_{ij}(t) = p_{ij} \forall t$ .

Likely, the simplest possible method to estimate the transition probabilities is inequality restricted least squares. Madansky (1959) notes that proportional data are inherently heteroskedastic so that weighted inequality restricted least squares (WLS) is appropriate. The inequality stems from the fact that the estimated transition probabilities are nonnegative. The WLS model is

$$\min Q(\mathbf{p}) = \sum_t \sum_j w_j \tilde{\epsilon}_j^2(t), \quad (2)$$

$$\text{s.t. } c_j(t) = \sum_i p_{ij} c_i(t-1) + \tilde{\epsilon}_j(t), \quad \forall j, t, \quad (3)$$

$$\begin{aligned} \sum_t \tilde{\epsilon}_j(t) &= 0, \quad \forall j \\ \sum_j p_{ij} &= 1, \quad \forall i, \\ p_{ij} &\geq 0, \quad \forall i, j \end{aligned} \quad (4)$$

$$w_j = \begin{cases} \frac{(T-r)}{\tilde{\epsilon}_j^2(t)} \\ \frac{T}{\sum_t c_j(t)} \\ \left[ \left( \frac{\sum_t c_j(t)}{T} \right) \left( \frac{1 - \sum_t c_j(t)}{T} \right) \right]^{-1} \end{cases} \quad (5)$$

Inequality restricted least squares is a special case of minimizing the objective function (2) subject to (3)–(5) when the weights,  $w_j$ , equal 1 for all classes. The equations in (3) are the Markov relation for each class and week, while the equations in (4) ensure that the residuals are mean zero, that all crop conditions are being modeled, and that the transition probability estimates are nonnegative. The specific weights considered in the estimation are presented in the equations in (5). As shown in (5), weights equal the inverse of the variance, inverse of mean sample proportions, and the inverse of the product of the mean sample proportions are all possible although the inverse of the variance is the most appealing. In the variance equation,  $\tilde{\epsilon}_j^2(t)$  are the (unweighted) inequality restricted least squares residuals found by minimizing (2) subject to (3) and (4) with the weights equal to 1 for each class. Therefore, a two-stage estimation is required with the residual variance from the first stage used in the weighting scheme for the second stage.

It should be noted that generalized least squares (GLS) would likely be a better choice as WLS does not explicitly account for serial correlation in the residuals. GLS is only possible if a state (i.e., an equation) is dropped from the system (2) subject to (3) and (4) due to the variance/covariance matrix of residuals being singular and therefore not invertible. However, a GLS approach also requires nonzero data for each class, and in a great many instances, zero data are reported for the “poor” and especially “very poor” classes. It is possible to substitute a small value for the zeros, say 0.001, but that approach was not taken here owing to the great number of substitutions that would be required.

Consequently, the WLS model above was used in the empirical application.<sup>3</sup>

As an alternative, we specify an information theoretic econometric model for the estimation of the transition probabilities that offers slightly more flexibility than the previous model. The maximum entropy (ME) model estimation is

$$\max H(\mathbf{p}, \mathbf{q}) - \sum_i \sum_j p_{ij} \ln p_{ij} - \sum_t \sum_j \sum_k q_{jk}(t) \ln q_{jk}(t), \quad (6)$$

$$\text{s.t. } c_j(t) = \sum_i p_{ij} c_i(t-1) + \sum_k u_k q_{jk}(t), \quad \forall j, t, \quad (7)$$

$$\begin{aligned} \sum_t \sum_k u_k q_{jk}(t) &= 0, \quad \forall j \\ \sum_j p_{ij} &= 1, \quad \forall i, \quad \sum_k q_k = 1, \quad (8) \\ p_{ij} &\geq 0, \quad \forall i, j \quad q_k \geq 0 \quad \forall k \end{aligned}$$

Here, the objective function in (6) is the maximization of entropy through the selection of two sets of probabilities: the  $p_{ij}$  as before, but also  $q_{jk}(t)$ , the probabilities associated with the error terms on the Markov relation in (7). As shown, the error term from (3) has been decomposed into the sum-product of a vector of error support values,  $u_k$ , and estimated probabilities,  $q_{jk}(t)$ .

Golan (2018) provides motivation for using the information concept of ME in econometric estimation and refers generally to the specification as info-metrics. Because the probabilities to be estimated are bounded between zero and one, the errors support values are bounded between  $\pm 1$  and the error support vector  $[-1 \ -\frac{1}{2} \ 0 \ \frac{1}{2} \ 1]$  is used. Lastly, the mean zero error, normalization, and nonnegativity equations in (8) are consistent with those in (4). The additional flexibility inherent in the system (6) subject to (7) and (8) is discussed below in the empirical application section of the paper. Lastly, the maximization of (6) subject to (7) and (8) does not admit a closed form solution for the estimator but is relatively easy to implement numerically.<sup>4</sup>

It should also be noted that the maximization of (6) subject to (7) and (8) is fundamentally different than squared error approaches such as the minimization of (2) subject to (3)–(5). In (6), the principle of ME involves finding the discrete probability distributions making up the rows of the transition probability matrix that obey the Markov relations but are otherwise as close as possible to uniform distributions. The proximity or divergence between the estimated and uniform distributions, measured according to (6), is not a true distance measure as

in (2) due to the former's violation of the triangle inequality.

### 3 | EMPIRICAL RESULTS

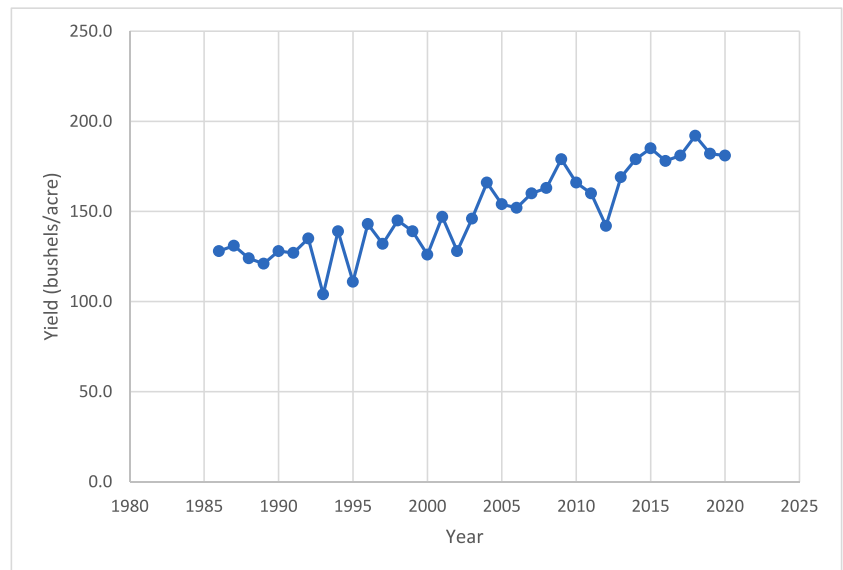
After estimating the transition probabilities describing the dynamics of crop condition, intrayear and final crop conditions can be estimated by forecasting the evolution of crop conditions using the estimates. Because crop conditions are released weekly during the growing season, weekly estimates of harvest time yield are possible by using the estimated transition probability matrix and a model describing how, for example, the final forecasted crop conditions translate to yield. To this end, an econometric model is specified wherein the final proportionate “excellent” and “good” crop conditions along with a trend variable are used to forecast the final yield.

#### 3.1 | Estimation of yields

As shown in Figure 1, since 1986, when crop condition data became available, average annual Nebraska corn yields appear to have a strong linear trend. Data from 2021 are not a part of the final estimation so they can be used to demonstrate model use and comparison with actual intrayear and final USDA yield estimates. The regression results presented in Table 1 show that, consistent with Schnitkey (2017), Irwin and Good (2017a, 2017b), and Irwin and Hubbs (2018) for corn and soybeans grown in Illinois, there is a statistically significant relationship between the final Nebraska corn yield and the late season proportion of the crop considered “excellent” and “good.” Panel A of Table 1 shows the results of an ordinary least squares (OLS) regression of final yield as a function of a trend variable and the proportion of the corn crop in Nebraska considered “excellent” and “good” during the last week in a given growing season for which conditions are provided. As shown, the variables are all highly statistically significant and the model explains much of the variation in yields.

Even so, the residuals are serially correlated as indicated by the Durbin–Watson test statistic. Panel B of Table 1 shows the regression results after correcting for the serial correlation using the Prais–Winsten approach. The Ljung–Box test is a test of the null hypothesis that all autocorrelation coefficients out to  $m = \ln N$  lags are zero indicating a white noise time series. As shown, the OLS model (Panel A) suggests a rejection of the null in favor of at least one lag having a nonzero autocorrelation coefficient while the corrected model in Panel B shows a failure to reject the null. The augmented Dickey–Fuller test

**FIGURE 1** Nebraska average corn yield (bushels per acre), 1986–2020.



**TABLE 1** Nebraska average corn yield (bushels per acre) as a function of time trend and final “excellent” and “good” crop condition percentages, 1986–2020.

<b>Panel A: OLS estimation</b>			
Variable	Estimate	Standard error	t-statistic
Intercept	75.639***	5.729	13.20
Trend	2.055***	0.104	19.70
Excellent	69.000***	10.861	6.34
Good	48.511***	10.449	4.64
Adjusted $R^2$	0.932	Ljung–Box	4.034**
$F_{3,31}$	157.000***	ADF (no intercept)	−4.167***
Durbin–Watson	1.342**	ADF (intercept)	−4.105***
Jarque–Bera	2.427	ADF (intercept + trend)	−4.045***
<b>Panel B: Prais–Winsten serial correlation estimation</b>			
Variable	Estimate	Standard error	z-statistic
Rho	0.334**	0.162	2.07
Intercept	74.290***	5.713	13.00
Trend	2.067***	0.140	14.79
Excellent	69.974***	9.320	7.51
Good	50.396***	9.817	5.13
Ljung–Box	0.087	ADF (no intercept)	−4.142***
Durbin–Watson	1.908	ADF (intercept)	−4.081***
Jarque–Bera	2.159	ADF (intercept + trend)	−4.024***

Abbreviation: ADF, augmented Dickey–Fuller.

\*\*Significance at the 5% level.

\*\*\*Significance at the 1% level.

statistics indicate that the residuals are stationary and the results of the Jarque–Bera test indicate there is no statistical support for nonnormally distributed errors.

The regression results in Table 1, Panel B, suggest that corn yield in Nebraska tends to trend up about

2 bushels per acre annually, likely due to better technology (inputs like seed and herbicides as well as better production practices). In addition, each 1% of final “excellent” (“good”) corn acreage in Nebraska results in a final estimated +0.7 (+0.5) bushels of corn per acre. To

provide some perspective, since 1986, the final “excellent” (“good”) proportion of corn acreage in Nebraska has averaged 19.8% (48.5%) with a range of 2–45% (24–69%). As an example, a high of 45% “excellent” corn crop condition at the end of the growing season relative to the sample average results in  $(45\% - 19.8\%) \times 69.974 = 17.6$  more bushels of corn yield per acre on average.

### 3.2 | Estimation of transition probabilities

Shown in Table 2 are the WLS (Panel A) and ME (Panel B) estimated transition probabilities for corn crop conditions in Nebraska. Given weekly crop condition reports, transition probabilities to distant states were precluded via constraints on the estimation. For example, “excellent” corn acreage in Nebraska cannot transition to the “fair” state in one week’s time although it can transition to the “good” state. Similarly, “good” corn acreage can transition to the “excellent” or “fair” states (in addition to remaining “good”) but cannot transition to the “poor” or “very poor” states in one week’s time. In either case, the matrices are diagonally dominant implying that crop conditions, on average, tend to stay the same week to week. Further, both estimations are consistent with crop

conditions being more likely to improve rather than deteriorate. Irrigation may play a significant role in this result because more of Nebraska’s corn crop is irrigated than in any other state.

As noted above, statistical testing for time homogeneity showed no evidence that a Markov model characterizing Nebraska corn crop condition is not a stationary Markov chain. This result is likely attributable to the extent of irrigation in Nebraska. Both interyear and intrayear testing were conducted using methods suggested by Kelton and Kelton (1984) and consisted of (a) estimating transition probabilities for the sample periods 1986–2003 and 2004–2021, (b) estimating annual matrices for the first and second half of each year (i.e., within a growing season), and (c) pooling all data and estimating monthly matrices for comparison.

Regarding (a) above, there is no statistical evidence in favor of rejecting the null hypothesis of a time homogeneous Markov chain. Estimating restricted (whole sample) and unrestricted (two subperiod samples) sum of squared errors results in  $F_{v_1, v_2} = 0.60275$  for  $v_1 = 20$  and  $v_2 = 2,840$  degrees of freedom yielding a  $p$ -value of 0.9137. This result suggests that the evolution of Nebraska corn conditions from the time when the survey began until the early 2000s is not statistically different than it is more recently. Similarly, for (b) above,

TABLE 2 Weekly stationary Nebraska corn weighted least squares (WLS) and maximum entropy (ME) crop condition transition probability estimates, goodness-of-fit statistics, and limiting distributions, 1986–2020.

Panel A: WLS estimates					
	Excellent	Good	Fair	Poor	Very poor
Excellent	0.9015	0.0985	0.0000	0.0000	0.0000
Good	0.0400	0.9224	0.0376	0.0000	0.0000
Fair	0.0000	0.0857	0.8509	0.0633	0.0000
Poor	0.0000	0.0000	0.1652	0.7566	0.0782
Very poor	0.0000	0.0000	0.0000	0.0927	0.9073
$R^2$	0.8359	0.8377	0.8056	0.8538	0.9534
$+\infty$	0.1886	0.4643	0.2034	0.0780	0.0657
Panel B: ME estimates					
	Excellent	Good	Fair	Poor	Very poor
Excellent	0.8619	0.1381	0.0000	0.0000	0.0000
Good	0.0527	0.8985	0.0487	0.0000	0.0000
Fair	0.0000	0.1167	0.7855	0.0978	0.0000
Poor	0.0000	0.0000	0.2932	0.5355	0.1713
Very poor	0.0000	0.0000	0.0000	0.2753	0.7247
$R^2$	0.8315	0.8296	0.7956	0.8402	0.9305
$+\infty$	0.1884	0.4938	0.2062	0.0688	0.0428

rejection of the null of a time homogeneous Markov chain only occurred in 2016 (out of 35 estimated annual matrices), but at a level of significance that is a little higher than the conventional 5% ( $p$ -value = 6.97%). This result suggests that each year, there is no statistical difference between early and late growing season transition probabilities. Lastly, for (c) above, the null hypothesis of a time homogeneous Markov chain was only rejected when comparing the June to July matrices, but again at a level of significance a little higher than 5% ( $p$ -value = 6.03%). This result suggests that, for example, the transition probability matrix describing corn conditions in Nebraska in any month is not statistically different than the next month. One might argue that comparing the weekly transitions in June to those in July does result in some difference, but the  $p$ -value is not overly compelling, and the null cannot be rejected in any of the other months. For this reason and the strength of the other tests, our conclusion is that corn crop conditions in Nebraska are a stationary Markov chain.

Also shown in Table 2 are the goodness of fit  $R^2$  values by class as well as the invariant distribution implied by each matrix estimation. The overall  $R^2$  for the WLS (ME) estimation was 0.8335 (0.8296). The invariant distribution is suggestive of the distribution of the underlying stochastic process given enough time passes. While the sort of time passage required makes little sense given that growing seasons ultimately come to an end long before the invariant distribution is obtained, the underlying dynamics do suggest, at least mathematically, the specific invariant distribution shown. This is possibly where the ME formulation may offer superior flexibility for estimation in that the invariant distribution information can be incorporated directly into the estimation.

To see this, consider that the invariant distribution is found by solving the Kolmogorov backward equations:  $\mathbf{P}' = \mathbf{\Lambda P}$  or forward equations:  $\mathbf{P}' = \mathbf{P A}$  with  $\mathbf{P}(0) = \mathbf{I}$  where  $\mathbf{P}$  is the estimated transition probability matrix,  $\mathbf{\Lambda}$  is an infinitesimal generator matrix, and  $\mathbf{I}$  is an identity matrix. In either case, the invariant distribution results in a matrix with identical rows equal to those shown in Table 2. In the ME model (6) subject to (7) and (8), additional constraints can be added so that the invariant distribution implied by a given transition probability matrix is consistent with exogenous long-term average proportions. We specify the sample averages for each class as the long-term average proportions and condition the estimation of the transition probabilities on this information in an iterative fashion. As a technical aside, the long-term averages are only obtainable with error, and additional error supports and error probabilities by class

similar to those in (6)–(8) are added to the ME model (6) subject to (7) and (8). The long-term averages used are as follows: 17.2% (excellent), 53.7% (good), 20.9% (fair), 5.6% (poor), and 2.7% (very poor).

A convenient way to compare the WLS and ME invariant distributions implied by the transition probability estimates to the long-term average distribution noted above is by calculating the Kullback–Leibler (KL) divergence measure. The closer the measure is to zero, the closer two distributions are to one another in terms of probability. In fact, two identical distributions have KL divergence equal to zero. Comparing the WLS (ME) invariant distribution to the long-term distribution results in a KL divergence equal to 0.0302 (0.0081). Therefore, the ME estimation results in a transition probability matrix that, in addition to capturing the intrayear dynamics, better accommodates the empirically observed long-term average proportions.

Lastly, using the estimated matrices from Table 2 and the actual 2021 Nebraska corn condition proportions as they occurred over Weeks 21–42 in 2021 and shown in Figure 2, forecasts of the final corn condition proportions can be estimated as crop condition information arrives. These forecasted final proportions can then be used in the regression equation from Table 1 to estimate weekly crop yields. Shown in Figure 3 are the results for both WLS and ME estimations as well as USDA forecasts of the final Nebraska corn yield which began in Week 30 and ended in Week 42 in 2021. The final value shown and labeled “FINAL” is not a forecast but rather the actual USDA estimated total bushels of Nebraska corn production in 2021 divided by the total acres planted in 2021.

As shown, Weeks 21 to 29 suggest both yield forecasts are around 189 bushels per acre and that the ME estimate is slightly more conservative than the WLS estimate (0.61 bushel per acre on average). Beginning in Week 30 and continuing until Week 37, the estimates are decidedly closer to one another averaging only 0.21 bushels per acre difference with the WLS estimate being the slightly more conservative of the two. In addition, both estimates are trending downward over this period, and this result is coincident with a slight decrease (increase) in the percent of Nebraska corn considered “good” (“fair”) (see Figure 2). From Weeks 37 to 42, there is virtually no difference between the two estimates (0.03 bushels per acre difference on average), and both are increasing likely due to the increase in the percent “excellent” category (see Figure 2) during those weeks.

In addition, the official USDA estimate for August (186 bushels per acre) was initially below those suggested

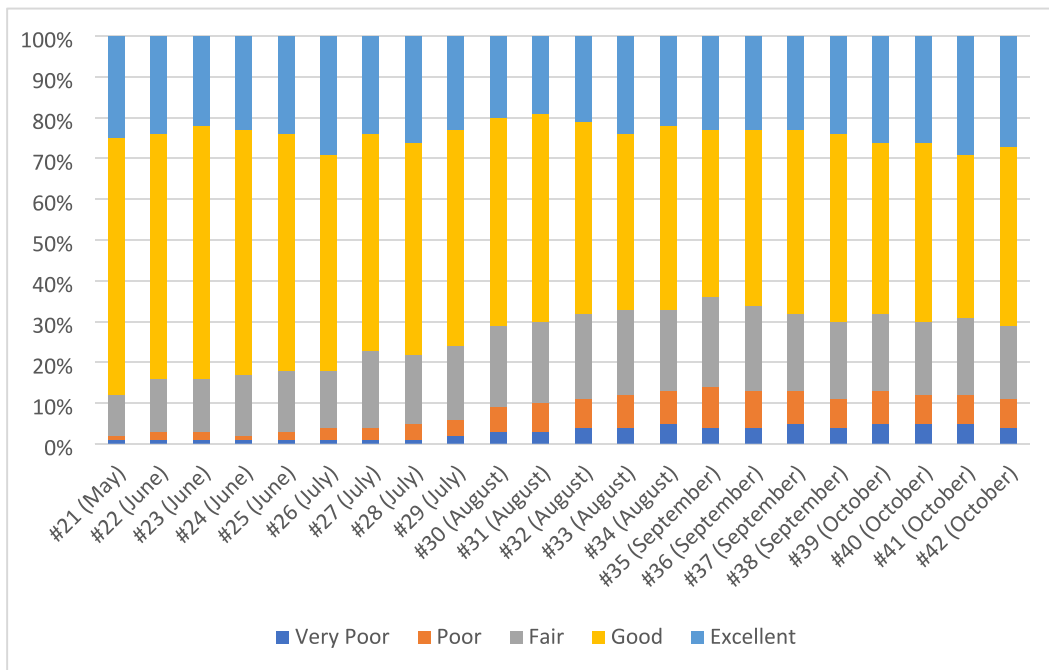


FIGURE 2 Weekly Nebraska crop conditions, 2021.

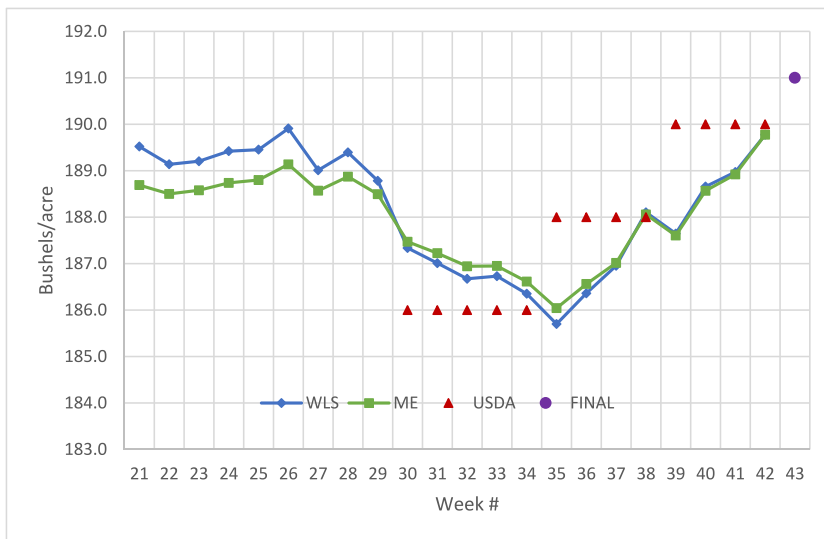


FIGURE 3 Weekly US Department of Agriculture and Markov chain estimates of final Nebraska corn yield (bushels per acre), 2021. ME, maximum entropy; WLS, weighted least squares.

by the two forecasts by about 1.5 bushels per acre. By contrast, the September (188 bushels per acre) and October (190 bushels per acre) USDA estimates were both above those forecasted by the Markov chain by about 2 bushels per acre. However, in both cases, the Markov chain yield estimates increased in the subsequent weeks to approximate the USDA value. Lastly, the final yield, as reported by the USDA (191 bushels per acre), was about 1 bushel per acre more than the October USDA estimate as well as the final Markov model estimates.

#### 4 | CONCLUSIONS

The research presented here adds to the ways in which crop yields can be forecasted by leveraging the flow of intraseason crop condition information. We show how crop conditions can be modeled as a Markov chain and show how to link the flow of intraseason crop condition information to final crop yields. Analytically, two different approaches are taken to estimate the transition probabilities making up the Markov chain. While both methods perform adequately, the ME formulation has



the added advantage of flexibility, being able to accommodate invariant distribution information as part of the estimation.

An empirical application of the methods presented for Nebraska corn production is also presented. The statistical evidence suggests that the Markov chain for Nebraska crop conditions is time homogeneous and both approaches to the estimation of transition probabilities result in similar forecasts. Even so, the ME approach for estimating transition probabilities is shown to be more flexible than weighted least squares. A simple econometric model relating final Nebraska corn yields to a trend variable and the proportion of “excellent” and “good” corn acreage forecasted from the Markov model fits the data well and provides accurate approximations to actual monthly USDA yield forecasts.

Future research should be directed at yield forecasts for other crops as well as total US crop production. In addition, other types of estimators are possible and may be more appropriate for other types of crops. Lastly, the method outlined here for intrayear crop yield forecasting likely has other applications. For example, bank loan risk ratings are often modeled as a first-order Markov chain that could be used to determine an important input for forecasting loan loss allowance.

## ACKNOWLEDGEMENTS

The author acknowledges the diligence of Derek Bunn, Editor of the *Journal of Forecasting*.

## DATA AVAILABILITY STATEMENT

Not applicable.

## ORCID

J. R. Stokes  <https://orcid.org/0000-0003-3000-0741>

## ENDNOTES

- <sup>1</sup> We note that unobservable  $a_i(t)$  are actually aggregated from the county to the state level before reporting  $c_i(t)$ .
- <sup>2</sup> The dynamic relationship between unconditional probabilities,  $\pi_j(t+1)$  and  $\pi_i(t)$ , is  $\pi_j(t+1) = \sum_i \pi_i(t) p_{ij}$ .
- <sup>3</sup> In matrix terms, the minimization of (2) subject to (3)–(5) results in the WLS estimator:  $\mathbf{P}_{WLS} = [\mathbf{c}'(t-1)\mathbf{H}'\mathbf{H}\mathbf{c}(t-1)]^{-1} + \mathbf{c}'(t-1)\mathbf{H}'\mathbf{H}\mathbf{c}(t)$  for generic weighting matrix  $\mathbf{H}$ .
- <sup>4</sup> It can be shown that the optimal transition probabilities are a function of the optimized Lagrange multipliers for the constraints in (7) and (8).

## REFERENCES

Anderson, T., & Goodman, L. (1957). Statistical inference about Markov chains. *The Annals of Mathematical Statistics*, 28(1), 89–110. <https://doi.org/10.1214/aoms/1177707039>

- Golan, A. (2018). *Foundations of info-metrics: Modeling, inference, and imperfect information*. Oxford University Press.
- Irwin, S., & Good, D. (2017a). When should we start paying attention to crop condition ratings for corn and soybeans? *farmdoc daily* (7):96, Department of Agricultural and Consumer Economics, University of Illinois at Urbana-Champaign May 24, 2017.
- Irwin, S., & Good, D. (2017b). How should we use within crop condition ratings for corn and soybeans? *farmdoc daily* (7):101, Department of Agricultural and Consumer Economics, University of Illinois at Urbana-Champaign June 1, 2017.
- Irwin, S., & Hubbs, T. (2018). What to make of high early season crop condition ratings for corn? *farmdoc daily* (8):108, Department of Agricultural and Consumer Economics, University of Illinois at Urbana-Champaign June 13, 2018.
- Kelton, D., & Kelton, C. (1984). Hypothesis tests for Markov process models estimated from aggregate frequency data. *Journal of the American Statistical Association*, 79(388), 922–928. <https://doi.org/10.1080/01621459.1984.10477112>
- Kruse, J., & Smith, D. (1994). Yield estimation throughout the growing season. Proceedings of the NCR-134 Conference on Applied Commodity Price Analysis, Forecasting, and Market Risk Management. Chicago, IL.
- Lee, T., Judge, G., & Zellner, A. (1970). *Estimating the parameters of the Markov probability model from aggregate time series data*. Amsterdam: North-Holland Publishing Co.
- Lehecka, G. (2014). The value of USDA crop progress and condition information: Reactions of corn and soybean futures markets. *Journal of Agricultural and Resource Economics*, 39(1), 88–105.
- Li, Y., & Andersson, J. (2020). A likelihood ratio and Markov chain-based method to evaluate density forecasting. *Journal of Forecasting*, 39(1), 47–55. <https://doi.org/10.1002/for.2604>
- Liu, C., Nassar, R., & Guo, M. (2015). A method of retail mortgage stress testing: Based on time-frame and magnitude analysis. *Journal of Forecasting*, 34(4), 261–274. <https://doi.org/10.1002/for.2326>
- Lo, C., Skindilias, K., & Karathanasopoulos, A. (2016). Forecasting latent volatility through a Markov chain approximation filter. *Journal of Forecasting*, 35(1), 54–69. <https://doi.org/10.1002/for.2364>
- Madansky, A. (1959). Least squares estimation in finite Markov processes. *Psychometrika*, 24, 137–144. <https://doi.org/10.1007/BF02289825>
- Matis, J., Birkett, T., & Boudreaux, D. (1989). An application of the Markov chain approach to forecasting cotton yields from surveys. *Agricultural Systems*, 29(4), 357–370. [https://doi.org/10.1016/0308-521X\(89\)90097-8](https://doi.org/10.1016/0308-521X(89)90097-8)
- Matis, J., Saito, T., Grant, W., Iwig, W., & Ritchie, J. (1985). A Markov chain approach to crop yield forecasting. *Agricultural Systems*, 18(3), 171–187. [https://doi.org/10.1016/0308-521X\(85\)90030-7](https://doi.org/10.1016/0308-521X(85)90030-7)
- Mattsson, A., & Thorburn, D. (1989). A simple check of the time homogeneity of Markov chains. *Journal of Forecasting*, 8(1), 65–72. <https://doi.org/10.1002/for.3980080106>
- Norwood, B., & Fackler, P. (1999). Forecasting crop yields and condition indices. Proceedings of the NCR-134 Conference on Applied Commodity Price Analysis, Forecasting, and Market Risk Management. Chicago, IL.

- Schnitkey, G. (2017). Still too early to say much about Illinois corn yields. *farmdoc daily* (7):95, Department of Agricultural and Consumer Economics, University of Illinois at Urbana-Champaign May 23, 2017.
- Tang, X., Hu, F., & Wang, P. (2018). Out-of-sample equity premium prediction: A scenario analysis approach. *Journal of Forecasting*, 37(5), 604–626. <https://doi.org/10.1002/for.2519>

## AUTHOR BIOGRAPHY

**J. R. Stokes** holds the rank of Professor at the University of Nebraska-Lincoln and conducts financial and

economic research related to banking, real estate finance, the measurement and management of credit risk, and applications of Markov chains.

**How to cite this article:** Stokes, J. R. (2023). A Markov chain model of crop conditions and intrayear crop yield forecasting. *Journal of Forecasting*, 1–10. <https://doi.org/10.1002/for.3052>