



United States Department of Agriculture

# Linking Public Data Sources to Create Localized Official Statistics

Greg Lyons & Dipak Subedi  
Economic Research Service  
October 24, 2018

The Findings and Conclusions in This Preliminary Presentation Have Not Been Formally Disseminated by the U. S. Department of Agriculture and Should Not Be Construed to Represent Any Agency Determination or Policy. This research was supported by the intramural research program of the U.S. Department of Agriculture, Economic Research Service.



# Motivation



**United States Department of Agriculture  
Economic Research Service**

ERS Home Topics Data Products Publications Newsroom Calendar Amber Waves Magazine

Home / Data Products / Farm Income and Wealth Statistics / Balance Sheet

United States	2014	2015	2016	2017	2018F
	\$1,000	\$1,000	\$1,000	\$1,000	\$1,000
<b>Farm sector debt</b>	345,201,354	356,738,041	374,164,212	393,048,069	406,854,605
Real estate	196,780,224	208,769,246	225,980,433	238,058,397	248,492,395
Commercial banks 1/	73,254,162	79,163,795	84,417,512	88,744,108	NA
Farm Credit System	88,797,518	96,662,553	103,749,537	107,653,783	NA
Farm Service Agency	4,325,689	4,857,770	5,914,514	6,054,097	NA
Farmer Mac	4,728,807	4,843,551	5,456,587	6,266,206	NA
Individuals and others 1/	12,517,927	9,956,273	12,494,207	13,463,931	NA
Storage facility loans	752,327	757,809	743,955	769,178	NA
Life insurance companies	12,403,795	12,527,497	13,204,121	15,107,093	NA
Nonreal estate	148,421,130	147,968,795	148,183,780	154,989,672	158,362,211
Commercial banks 1/	70,737,959	73,177,901	73,233,553	73,294,843	NA
Farm Credit System	47,887,186	48,283,041	49,376,260	51,180,555	NA
Farm Service Agency	3,550,210	3,748,543	3,783,890	3,958,398	NA
Individuals and others 1/	26,245,776	22,759,310	21,790,077	26,555,877	NA

**Footnotes**  
 Data as of August 30, 2018  
 F = Forecast values.  
 NA = Data are not available/applicable.  
 Values are rounded to the nearest thousand. When 'Real (2018 dollars)' is selected, nominal values are adjusted for inflation using the chain-type GDP deflator, base year=2018.  
 1/ Beginning with 2012 estimates, farm sector debt held by savings associations is reported with the commercial bank lender group instead of the individuals and others grouping.

[USDA/ERS Farm Income and Wealth Statistics](http://www.ers.usda.gov)

The Economic Research Service produces national balance sheets as part of our Farm Income and Wealth Statistics data products

Objective: find a method to procure state-level estimates through better use of existing reports and new disaggregation methods of administrative data

Focus: 85% of loan volumes held by Commercial/Savings Banks, the Farm Service Agency and the Farm Credit System



# Challenges for Top Institutional Lenders

Institution	Issue
Farm Service Agency	State-level data exists, but not in readily available format
Commercial/Savings Banks	Data is aggregated by bank, not by state. State-level data can be imputed with regulatory sources
Farm Credit System	Data is aggregated by bank, not by state. Limited regulatory information. State values must be estimated using other means (e.g. surveys)



# Data Sources

## **Commercial/Savings Banks**

Call Report Data: 1976 – 2018

- Federal Reserve Bank of Chicago; Federal Financial Institutions Examination Council

Summary of Deposits: 1994 - 2018

- Federal Deposit Insurance Corporation

Community Reinvestment Act: 1997 – 2018

- Federal Financial Institutions Examination Council

Home Mortgage Disclosure Act: 1999 – 2016

- Federal Financial Institutions Examination Council

## **Farm Service Agency**

Monthly Management Summary Reports: 2003 – 2018

- Farm Service Agency

## **Farm Credit System**

Call Report Data: 2005 – 2018

- Farm Credit Administration

## **Other Sources**

Census of Agriculture: 1992 – 2012

- USDA National Agricultural Statistical Service



# Accessing State-Level Information From PDFs



# Using R to Scrape PDFs from the Farm Service Agency

## Sample PDF

Objective: Read in data frame of tabular data in PDF of loan data for states

Tabular data in FSA PDFs were accessed by

- Transforming the PDFs into a data frame containing lines of text
- Indexing start and end of table using regular expressions
- Coercing fixed width data into columns

The image shows a sample PDF document containing a table of data. The table has several columns and rows, with data points that are difficult to read due to the low resolution and blurriness of the scan. The table appears to be a summary of loan data for various states, with columns likely representing state names, loan amounts, interest rates, and other financial metrics. The text is organized in a structured, tabular format typical of a data report.



# Using R to Scrape PDFs from the Farm Service Agency

## Example – PDF Scrape

Packages used:

Pdftools ← PDF to text

Stringr ← string manipulation

Full script had to account for quirks, such as

- Changes to table over time
- “West Virginia”

```
1 #Example: Pulling table of state data from PDF
2 #Convert PDF into textfile
3 textfile <- pdf_text(filepath)
4
5 #Creating new row in dataframe for each space
6 flatfile <- strsplit(textfile, "\n")
7
8 #Searching for page containing table of interest TableName
9 for(i in 1:length(flatfile)){
10
11 #Identifying start of table TableName
12 if(grep(sprintf(TableName), flatfile[[i]][3])){
13
14 #Looping through rows of interest using regular expressions
15 for(j in grep("ALABAMA", flatfile[[i]]):grep("WYOMING", flatfile[[i]])){
16
17 #Converting each each row into a vector
18 vector <- unlist(strsplit(str_replace(gsub("\\s+",
19 " ", str_trim(tolower(flatfile[[i]][j]))), "B", "b"), " "))
20
21 #Binding each row to a new dataframe, outfile
22 outfile <- rbind(outfile, vector, deparse.level = 0, stringsAsFactors=FALSE)
23 }
24 }
25 }
```



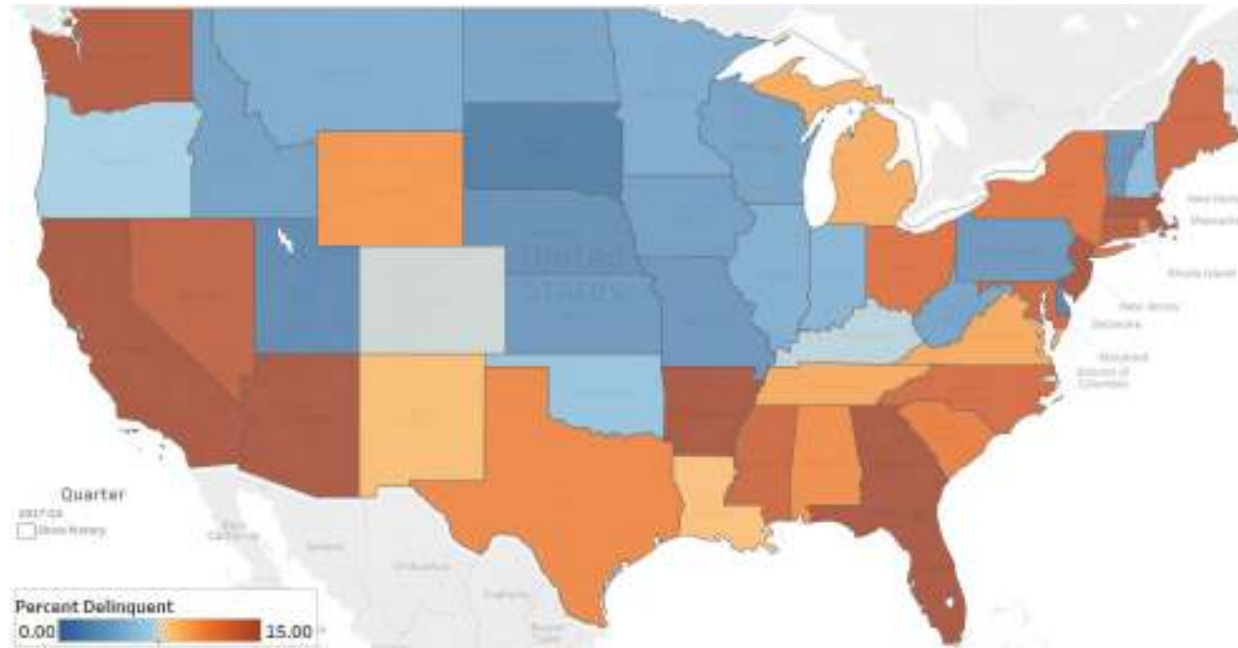


# Using R to Scrape PDFs from the Farm Service Agency

## Delinquency Rates for FSA Production Loans – Q2 2018

Benefits of this approach:

- Parameterized code allows for automatic quarterly updates
- Additional variables or tables can be extracted with minimal code changes
- No need for intermediate tables for data visualizations





# Disaggregating Bank-Level Data with Regulatory Information



# Using R to Disaggregate Commercial Bank Call Report Data

Call reports list information by headquarters, not where loans occur

Solution: disaggregate call reports information into counties using regulatory information that captures bank presence by county and re-aggregate at the state level

Overarching Process:

- 1) Read in data
- 2) Cleaning and imputation
- 3) Assigning county shares
- 4) Many-to-one merges
- 5) Calculation of county shares
- 6) Re-aggregation and upload

Original Data – Loan Volume by Institution



Disaggregated Data – Loan Volume by County



# Accessing Call Report/Regulatory Information

Most sources used are contained in zip files that have URLs that can be used for direct access

Packages used:

RCurl ← url access

SASxport ← read xport files

Section Process

- Download zip file to temporary directory
- Identify index using regular expressions
- Merge pulled schedules
- Delete temporary files

Significant cleaning, but uses simple methods

## Example – Zip File Extract

```
1 #Example - reading in call report data from multiple schedules
2 #Adding file.exists helps avoid failures
3 if(file.exists(call_report_url)){
4
5     #Create a temporary directory and download the full zipfile
6     td <- tempdir()
7     tf = tempfile(tmpdir=td, fileext=".zip")
8     download.file(call_report_url,tf,mode="wb")
9
10    #Find the index number for the name of the first schedule and read in the data
11    fname = unzip(tf, list=TRUE)$Name[grep("SCHEDULE A",unzip(tf, list=TRUE)$Name)]
12    unzip(tf, files=fname, exdir=td, overwrite=TRUE)
13    fpath = file.path(td, fname)
14    sched_a = read.xport(fpath)
15
16    #Repeat the process for the second schedule
17    fname = unzip(tf, list=TRUE)$Name[grep("SCHEDULE B",unzip(tf, list=TRUE)$Name)]
18    unzip(tf, files=fname, exdir=td, overwrite=TRUE)
19    fpath = file.path(td, fname)
20    sched_b = read.xport(fpath)
21
22    #Combine schedules
23    all_schedules <- Reduce(function(x, y) merge(x, y, all=TRUE), list(sched_a,sched_b))
24
25    #Delete all files in temporary directory
26    do.call(file.remove, list(list.files(td, full.names = TRUE)))
27 }
```



# Assigning County Proportions

Merging call report data with regulatory data allows us to proxy for an institution's regional lending

One source will contain aggregated information you are attempting to disaggregate, and the other(s) contain disaggregated information that can be used to proxy for regional dispersion

## Section Process

- Add column to regulatory table containing sum by group
- Merge in call report data (many to one)
- Multiply across for county share

## Largest Lender: Call Report Information



## Largest Lender: Regulatory Information



# Assigning County Proportions

Packages used:

Rodbc ← SQL server connection

Challenge: often requires strong assumptions, in-depth knowledge and significant cleaning before merge

Note: many-to-one merges can result in tables that are too large for individual machines to hold in memory

## Example – Assigning County Shares with Share of Total

```
1 #Example - using regulatory information to assign shares
2 #need to create sum of volume by unit_id and add as column
3 df$unit_sum <- ave(df$volume, df$unit_id, FUN=sum)
4 #Divide volume (county level) by this total for share
5 df$county_share <- df$volume / df$unit_sum
6 #checking to see if shares by unit add to 1
7 df$check <- ave(df$county_share, df$unit_id, FUN=sum)
8 #If correctly disaggregated, table should only include 1s
9 table(df$check)
10 #merge in call report information where reg info exists
11 df_with_cr <- merge(df, call_reports, by="unit_id", all.x=TRUE, all.y=FALSE)
12 #county-level volumes
13 df_with_cr$loan_volume_share <- df_with_cr$county_share * df_with_cr$loan_volume
14 #Saving database out to SQL - where dbhandle is database handle
15 sqlsave(dbhandle, df_with_cr, "disagg_call_reports", fast=TRUE, append=TRUE,
16         rownames=FALSE)
```





# Refining Administrative Data with Surveys





# Combining Administrative Data with Surveys

Similarities to commercial bank data:

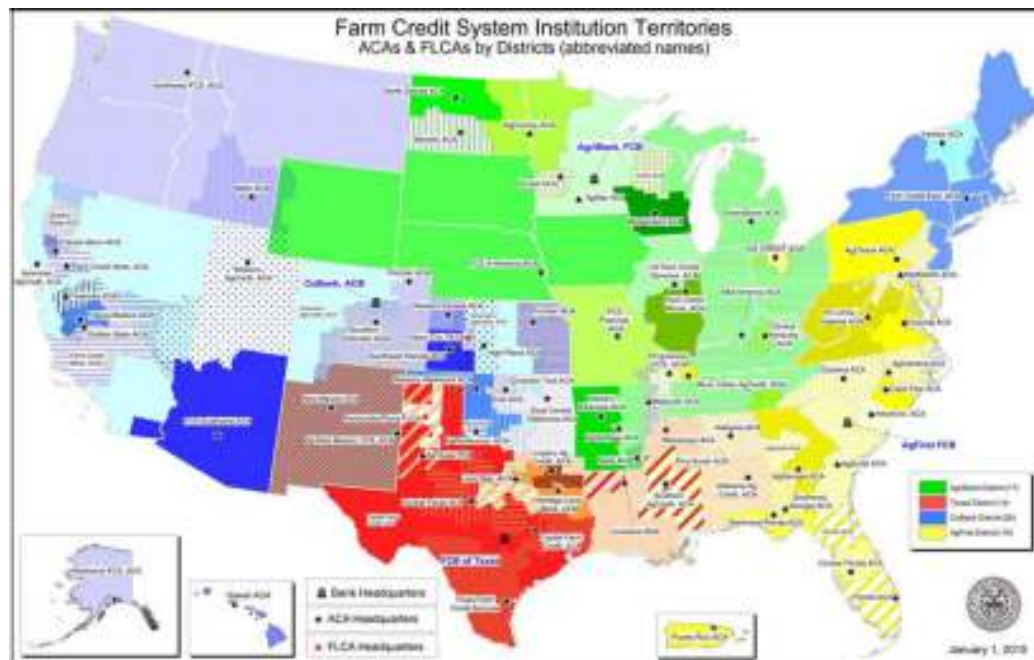
- 1) Data reported at the institution level
- 2) Same importation method

Differences:

- 1) Limited regulatory information

Overarching Process:

- 1) Read in data
- 2) Survey analysis
- 3) Assign state shares



# Using R to Disaggregate Commercial Bank Call Report Data

Packages used:

survey ← survey analysis

Useful specifically when survey total is less reliable than what is reported in administrative data, but produces valid proportions by group

Important to know survey limitations to know what mitigating factors to use (e.g. moving averages)

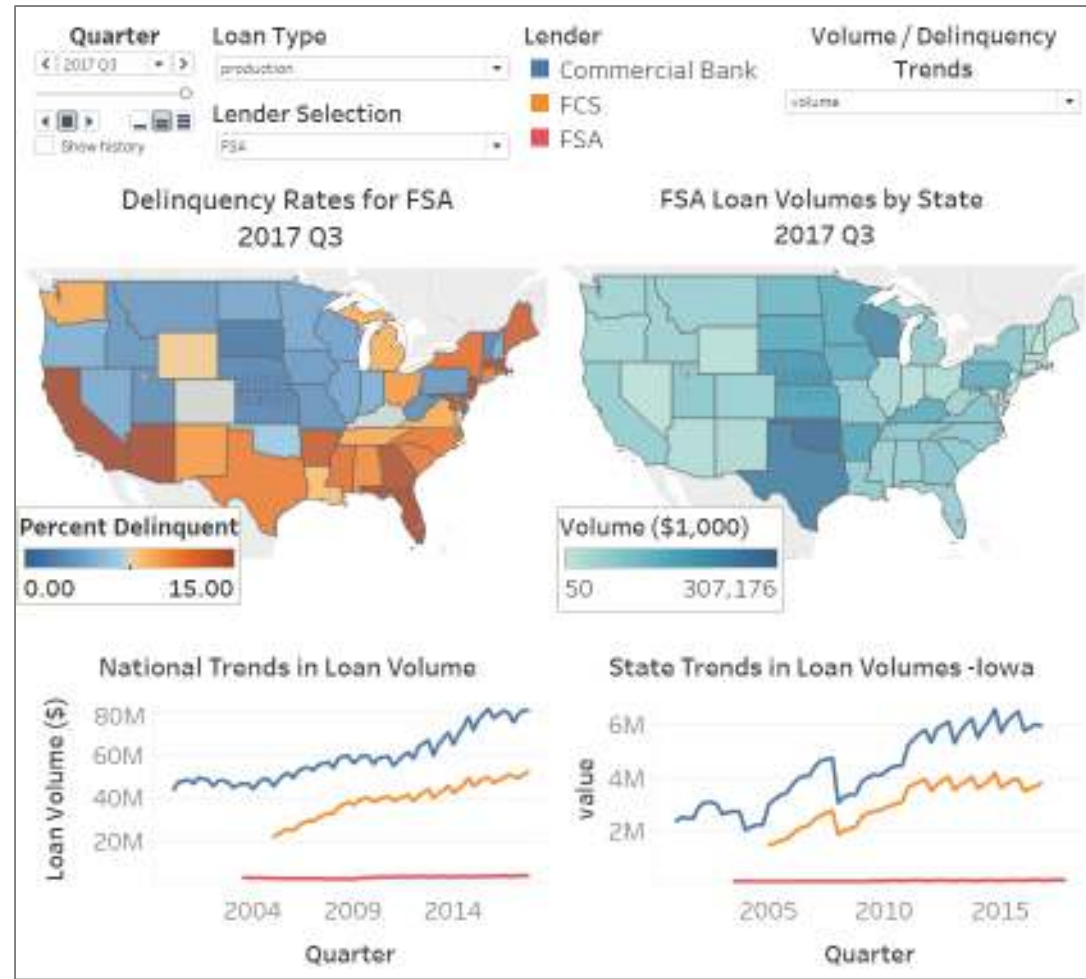
## Example – Assigning State Share using Survey Share

```
1  ### Creating state shares based on national total
2  # read in survey data
3  data <- read.csv(surveypath,header=T)
4  # create survey design
5  survey.design <- svydesign(id=~id,data=data,weights = ~weights)
6  # summing debt by state
7  state_debt <- svyby(~debt,~state,survey.design,svytotal)
8  # finding share of debt by state
9  state_debt$debt_share <- state_debt$debt / sum(state_debt$debt)
10 # applying state share of debt to administrative total
11 state_debt$admin_share <- state_debt$debt_share * national_total
```



# Combining All Methods

## Example: Dashboard Mockup



## Potential Use Cases

- Extension of ERS data products
- Creation of new ERS visualizations
- Use for broader research purposes



## Contact Information

Greg Lyons, ERS/USDA  
(202) 694-5147  
[greg.lyons@ers.usda.gov](mailto:greg.lyons@ers.usda.gov)



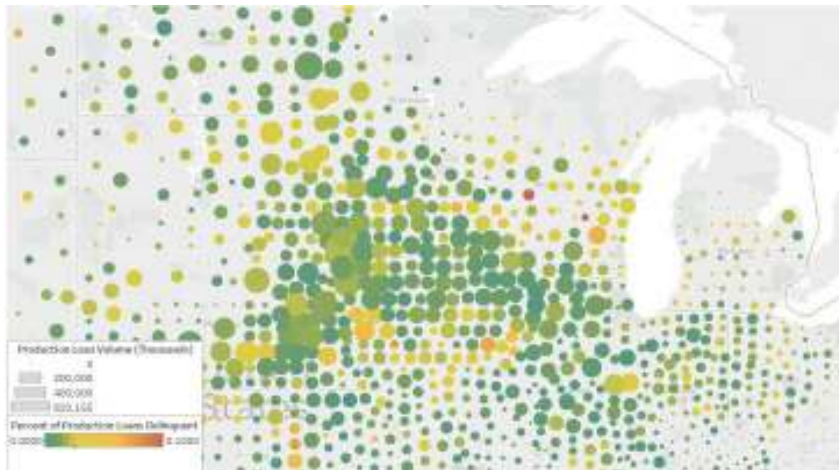


# Extensions: Financial Stress Factors

Real Estate Loan Delinquency Rate by County, Q2 2018



Production Loan Delinquency Rate by County, Q2 2018



Disaggregation methods can be extended to other schedules to look at regional financial stress, including:

- Delinquent loan volumes
- Loans in nonaccrual status
- Charge-offs



# Extensions: Bank Branch Closures

Number of Closed Branches by County, 2017



Subsequent research will pair county-level loan volumes and delinquencies with the bank branch closures from the FDIC's Reports of Structure Changes

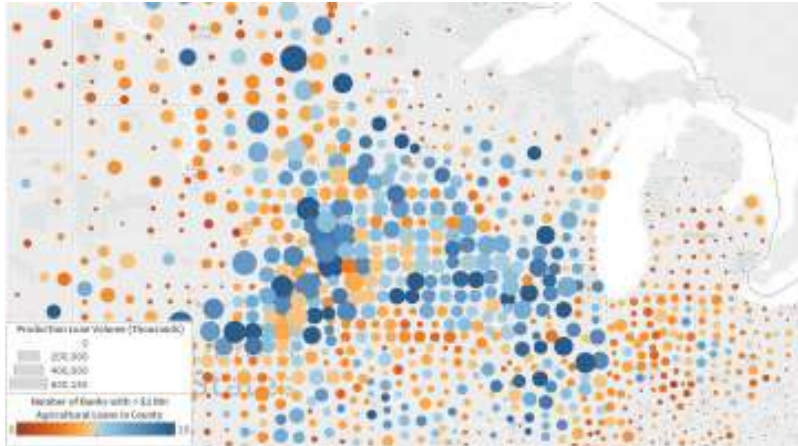
Primary aim is to understand the relationship between agricultural loan performance and bank branch closures



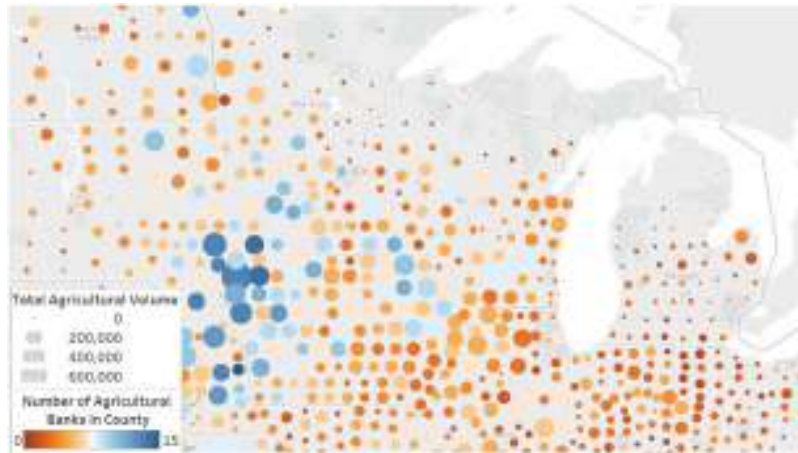


# Extensions: Competition in Lending

Number of Banks with > \$1Mn Loans in County



Number of FDIC-designated Ag. Banks in County

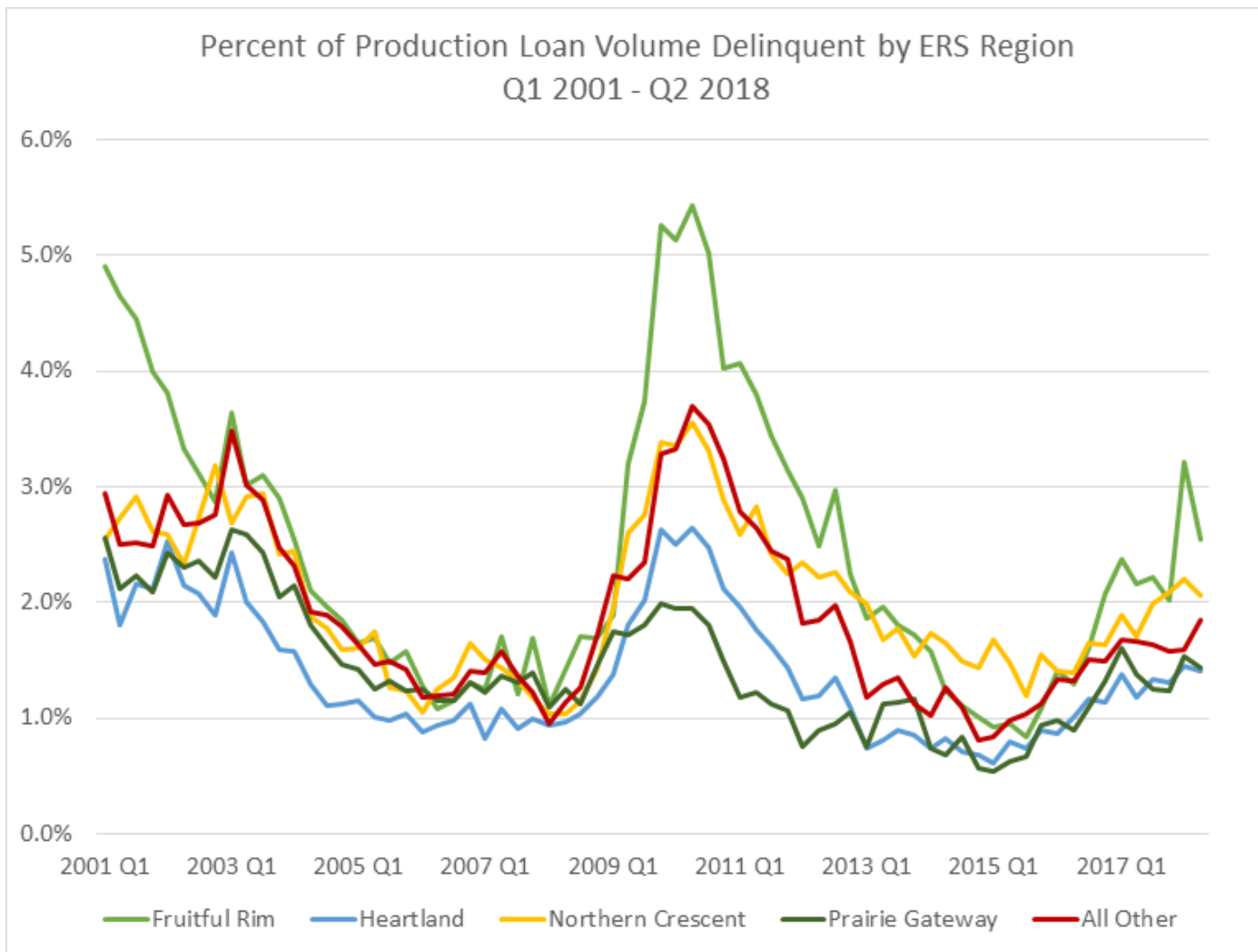


Can be used to check the robustness of agricultural credit markets across counties

- Number of institutions with x in agricultural lending
- Number of agricultural-focused lending institutions



# Appendix: Production Loan Delinquencies by ERS Production Region



# Appendix: Real Estate Loan Delinquencies by ERS Production Region

