



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 \$1,000	2015 \$1,000	2016 \$1,000	2017 \$1,000	2018F \$1,000
Farm sector debt	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

FILE: FPLFC (03/14) DATE: 05/03/18 PART 2

FSA - FARM LOAN PROGRAM - NATIONAL SUMMARY TOTAL REP DISTRICT LOAN BORROWER SUMMARY BY STATE BASED ON RCH40 AC OF: 04/30/18

STATE	UNPAID PRINCIPAL	UNPAID INTEREST	TOTAL P & I OWED	AMOUNT P & I OWED	DOLLAR DELIN. CNTY	PER-ROWING DELQ	NRB OF BORROWERS	NRB DELIN. CNTY	PER-ROWING DELQ
ALASKA	106,311,472	3,627,211	109,938,684	8,270,654	7,52	1,965	362	18,42	
ARIZONA	6,327,676	325,384	6,653,060	954,709	14,38	63	11	17,46	
ARKANSAS	37,131,946	1,521,340	38,653,286	6,098,457	10,79	395	103	26,08	
CALIFORNIA	276,217,957	17,651,474	293,871,431	39,194,888	13,48	2,177	629	29,08	
COLORADO	150,162,978	23,557,824	173,720,802	36,374,815	20,94	1,438	217	15,41	
CONNECTICUT	167,051,938	3,025,442	170,077,380	6,856,976	3,43	1,019	165	16,39	
DELAWARE	5,112,899	317,274	5,430,172	741,834	13,66	55	16	29,09	
FLORIDA	16,177,772	103,622	16,281,394	31,472	19	79	6	6,33	
GEORGIA	109,347,348	4,647,320	113,994,668	13,310,741	11,69	995	246	26,73	
ILLINOIS	169,625,757	8,039,806	177,665,563	20,517,632	12,13	1,478	356	24,09	
INDIANA	18,433,021	217,642	18,650,663	312,921	1,68	399	44	11,03	
IOWA	149,792,742	2,271,644	151,064,386	3,152,918	6,29	1,189	108	10,51	
KANSAS	371,569,677	4,125,783	375,695,460	2,371,155	63	2,813	102	3,63	
KENTUCKY	196,214,813	1,921,277	198,136,090	1,922,904	97	1,339	13	4,45	
LOUISIANA	683,981,863	6,888,556	690,870,418	8,210,856	1,19	4,976	307	6,17	
MAINE	576,956,435	7,746,035	584,702,470	6,343,837	13,05	3,822	235	6,13	
MARYLAND	408,351,144	9,108,598	417,459,742	17,880,307	4,28	4,375	735	16,80	
MASSACHUSETTS	99,536,056	6,448,413	105,984,469	14,057,913	13,91	1,542	287	18,61	
MICHIGAN	35,312,967	2,021,137	37,334,104	1,550,731	9,38	316	69	21,84	
MINNESOTA	76,149,236	715,309	76,864,545	1,792,639	6,67	361	34	11,32	
MISSISSIPPI	47,039,894	2,486,637	49,526,531	7,080,295	14,30	364	118	32,42	
MISSOURI	217,425,737	2,775,030	220,200,766	10,648,751	4,63	1,709	214	12,52	
MONTANA	491,807,648	7,893,325	500,000,973	10,939,661	2,19	3,086	397	12,67	
NEBRASKA	70,883,856	6,115,963	77,000,000	82,000,819	13,09	960	185	17,64	
NEVADA	334,844,141	4,517,993	339,362,134	2,843,669	84	2,733	157	5,74	
NEW HAMPSHIRE	144,923,531	5,355,147	150,278,678	7,940,869	1,28	987	109	11,04	
NEW JERSEY	642,431,742	7,276,496	649,708,238	9,188,368	1,41	4,210	348	8,27	
NEW MEXICO	23,732,750	938,489	24,671,239	1,743,951	7,19	233	28	14,97	
NEW YORK	15,954,637	344,331	16,298,968	788,327	1,79	387	20	10,71	
NORTH CAROLINA	25,965,275	2,941,886	28,907,161	4,700,555	16,78	214	69	21,78	
NORTH DAKOTA	77,877,510	2,136,998	80,014,508	1,842,833	4,80	722	154	15,79	
OHIO	355,984,371	7,843,180	363,827,551	17,997,325	10,02	1,172	306	26,11	
OREGON	139,667,865	4,654,861	144,322,726	16,769,374	11,66	1,166	305	26,16	
PENNSYLVANIA	202,464,133	4,637,887	207,102,020	8,466,114	3,32	1,128	198	12,88	
RHODE ISLAND	212,189,436	3,241,858	215,431,294	6,113,628	3,39	987	157	8,45	
SOUTH CAROLINA	1,010,631,186	21,775,656	1,032,406,842	26,149,844	2,13	6,884	956	13,16	
TENNESSEE	103,708,626	1,691,691	105,400,317	2,256,121	1,14	789	95	9,51	
TEXAS	1,074,286,965	4,863,448	1,079,150,413	8,093,385	2,86	1,700	244	14,35	
UTAH	162,841,826	74,680,113	237,521,939	120,669,701	50,91	1,745	1,199	68,71	
VIRGINIA	7,132,146	397,137	7,529,283	853,976	11,41	75	11	15,07	
WASHINGTON	138,771,874	5,499,204	144,271,078	14,313,793	10,16	1,066	219	21,77	
WEST VIRGINIA	484,922,793	6,716,245	491,639,038	4,241,396	38	2,065	188	6,47	
WISCONSIN	139,011,974	17,776,405	156,788,379	10,614,677	6,46	1,014	192	24,79	
WYOMING	139,011,974	17,776,405	156,788,379	10,614,677	6,46	1,014	192	24,79	
TOTAL	10,557,219,080	327,178,795	10,884,397,875	573,603,170	5,27	84,478	12,222	14,47	

SOURCE: FSA - ECOM FOCUS REPORT *RCH40* DATA BASE.
FOR INTERNAL DISTRIBUTION ONLY



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 }
```

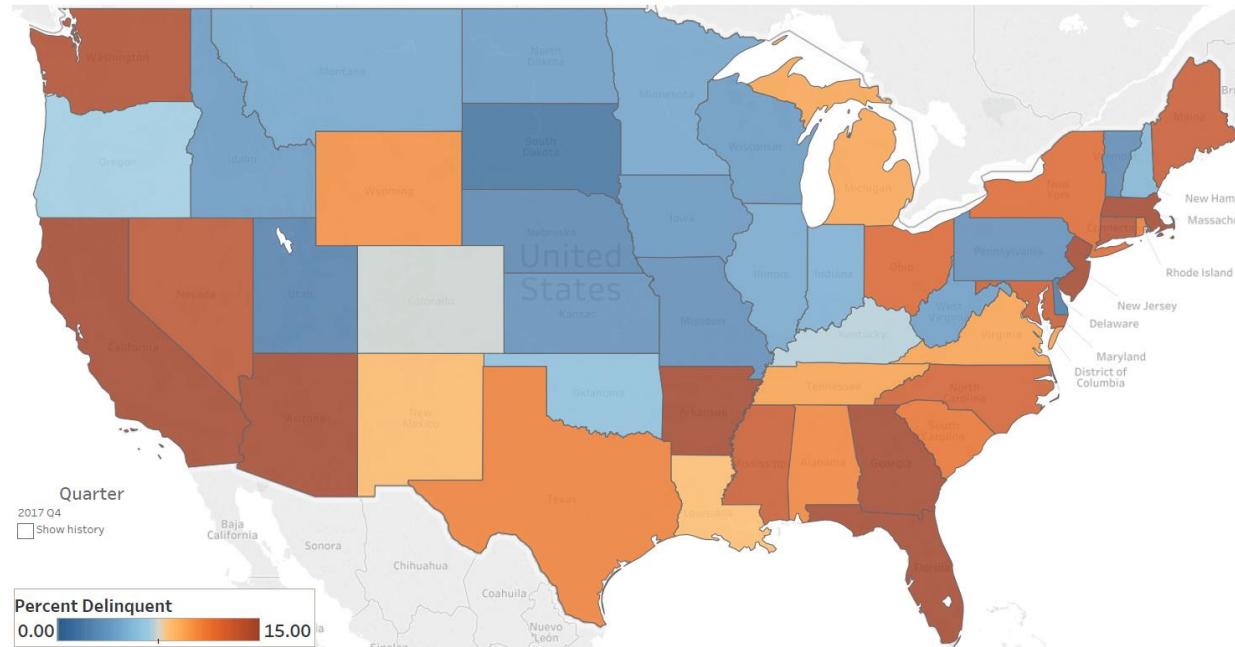


Using R to Scrape PDFs from the Farm Service Agency

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

Delinquency Rates for FSA Production Loans – Q2 2018



Disaggregating Bank-Level Data with Regulatory Information



Using R to Disaggregate Commercial Bank Call Report Data

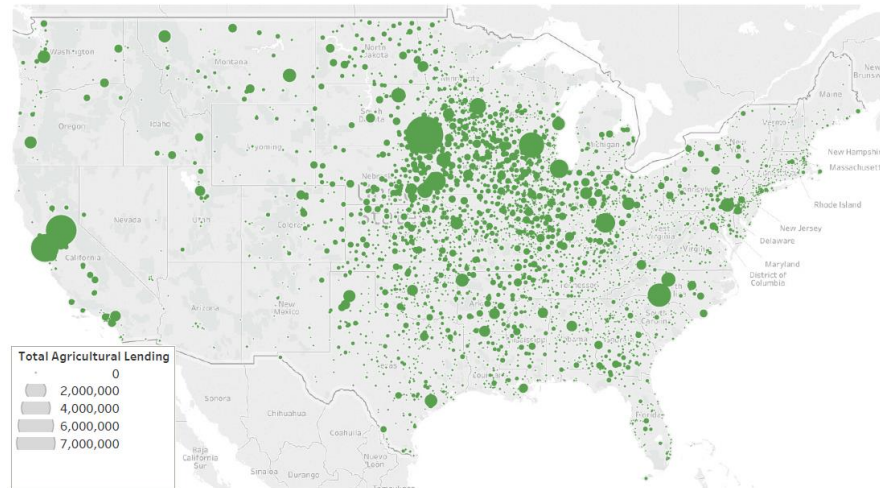
Call reports list information by headquarters, not where loans occur

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

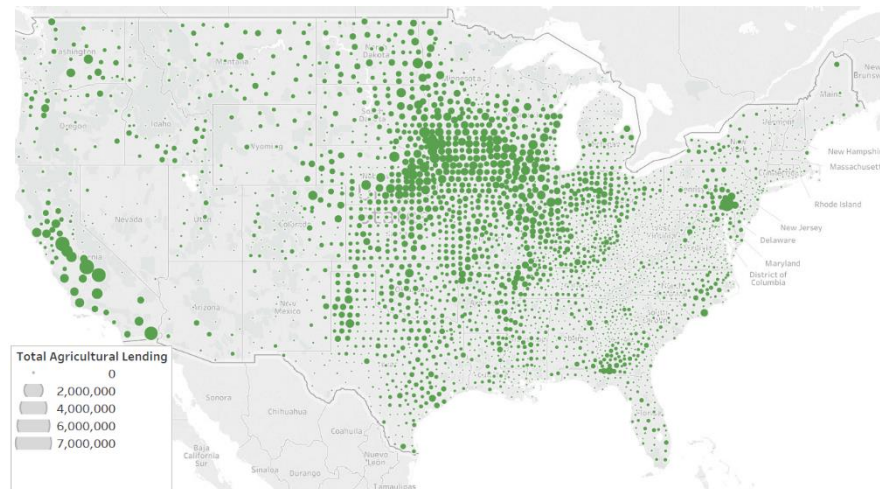
Overarching Process:

- 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

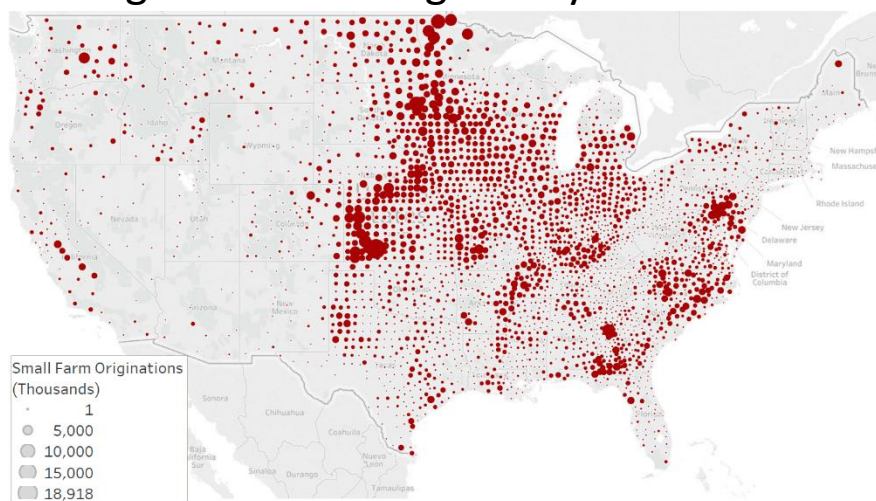
Section Process

- 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:

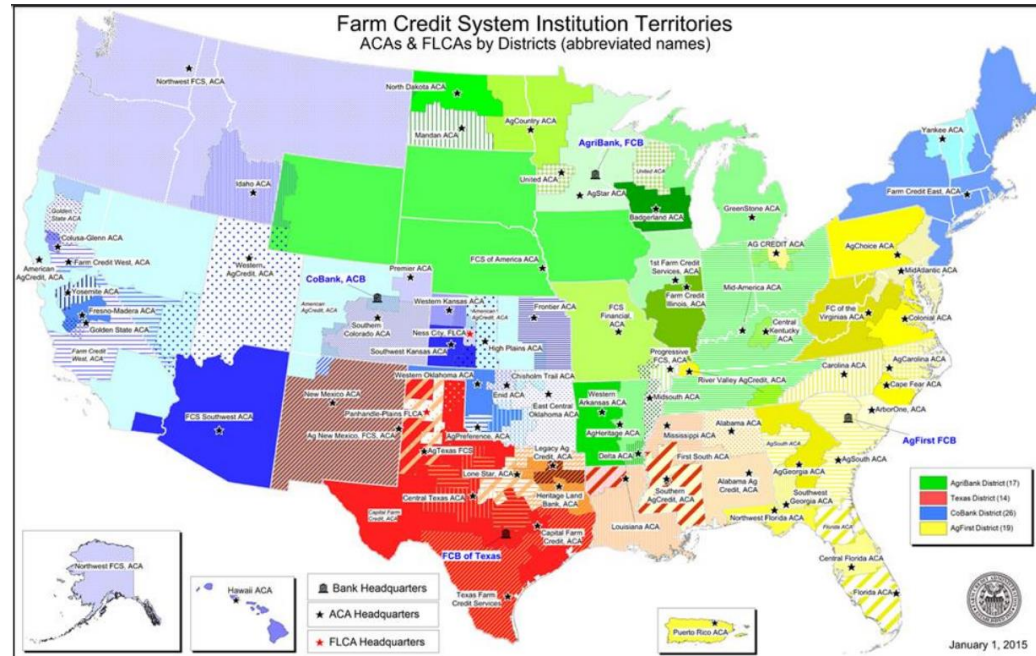
- 1) Data reported at the institution level
- 2) Same importation method

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)

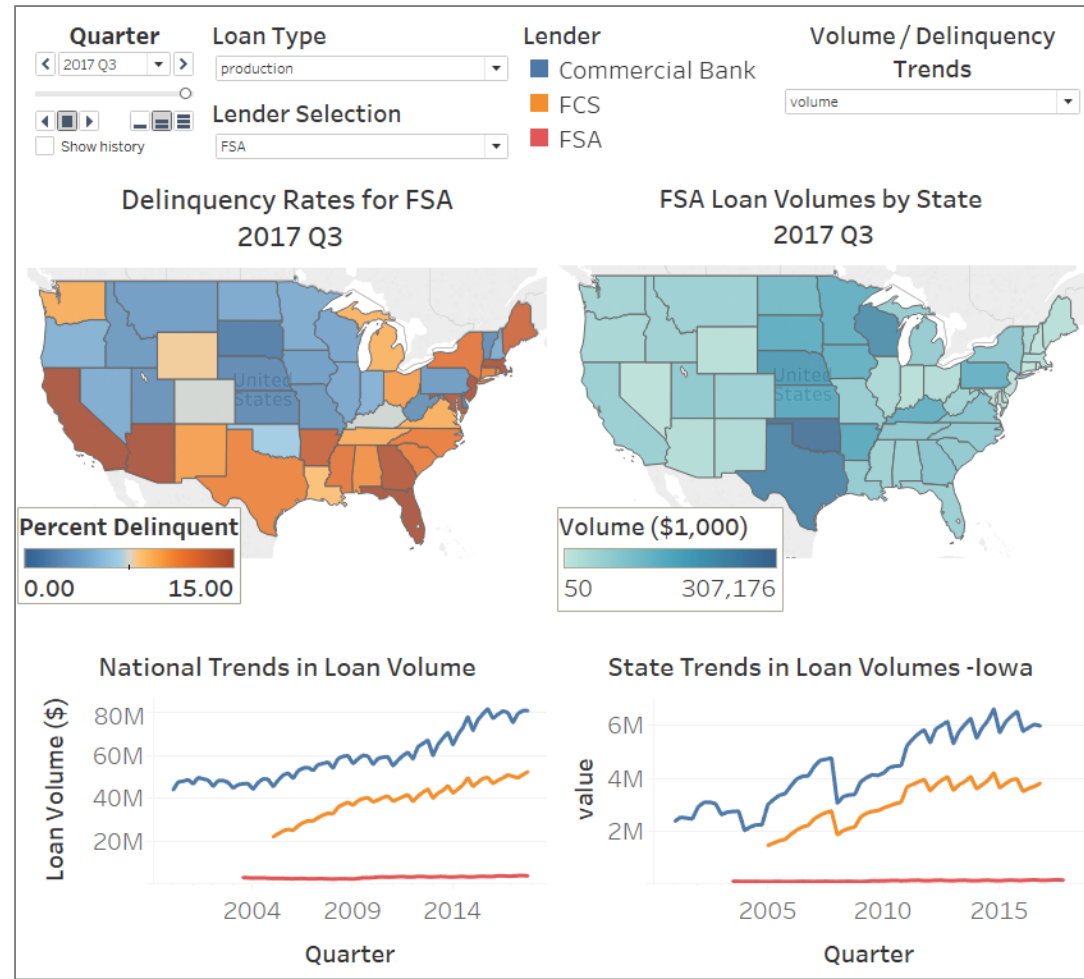
Example – Assigning State Share using Survey Share

```
1  ### Creating state shares based on national total
2  # read in survey data
3  data <- read.csv(survey_path,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



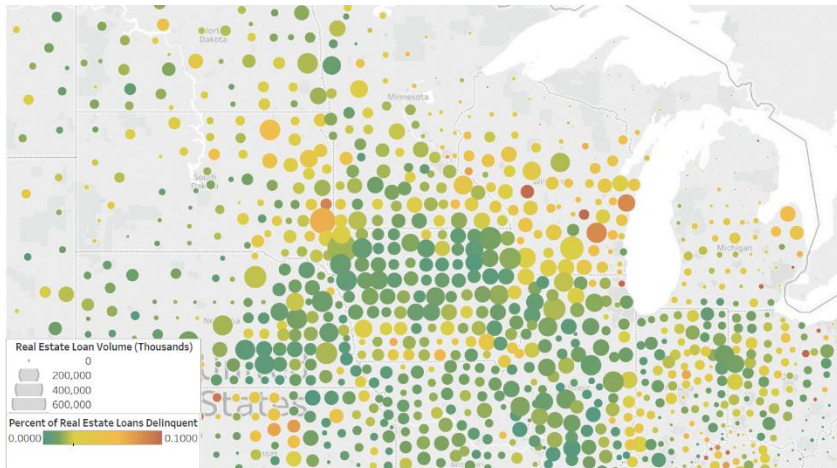
Contact Information

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

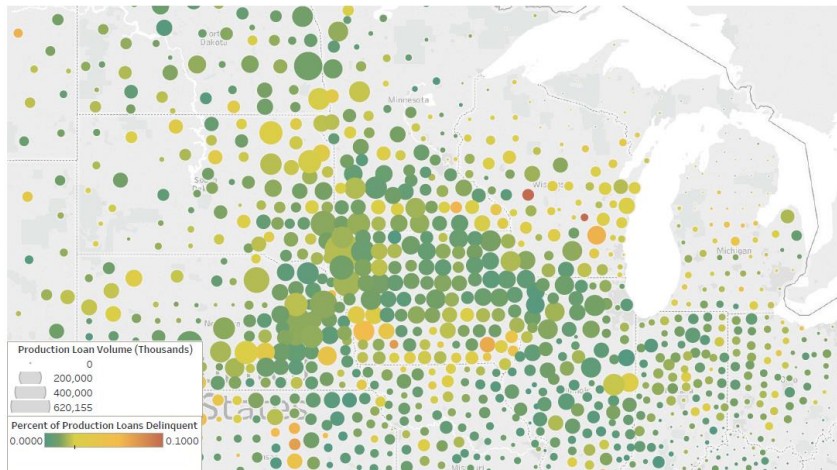


Extensions: Financial Stress Factors

Real Estate Loan Delinquency Rate by County, Q2 2018



Production Loan Delinquency Rate by County, Q2 2018



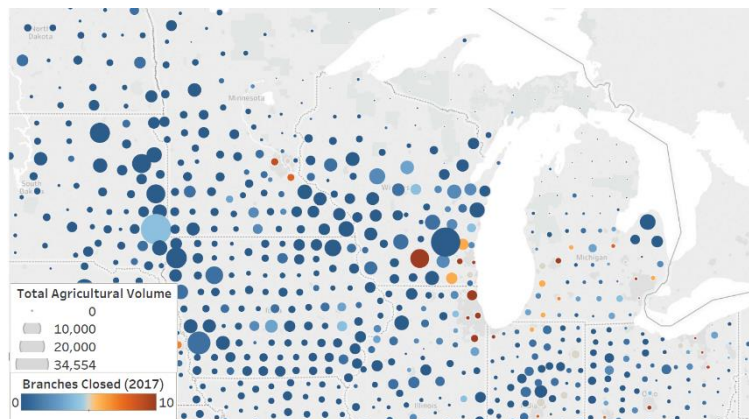
Disaggregation methods can be extended to other schedules to look at regional financial stress, including:

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



Extensions: Bank Branch Closures

Number of Closed Branches by County, 2017



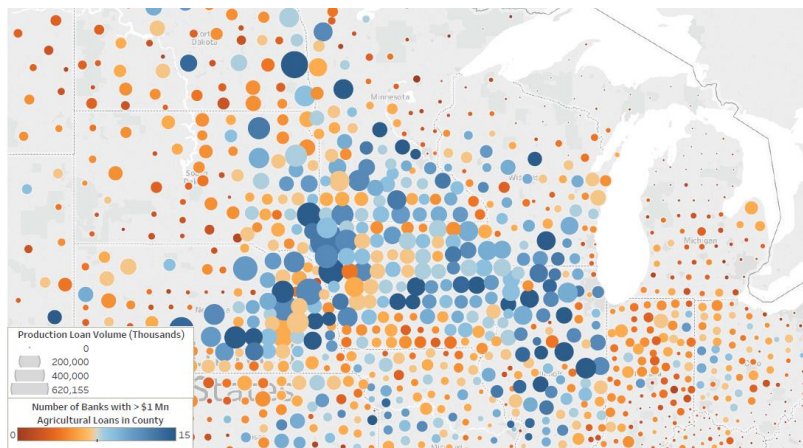
Subsequent research will pair county-level loan volumes and delinquencies with the bank branch closures from the FDIC's Reports of Structure Changes

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



Extensions: Competition in Lending

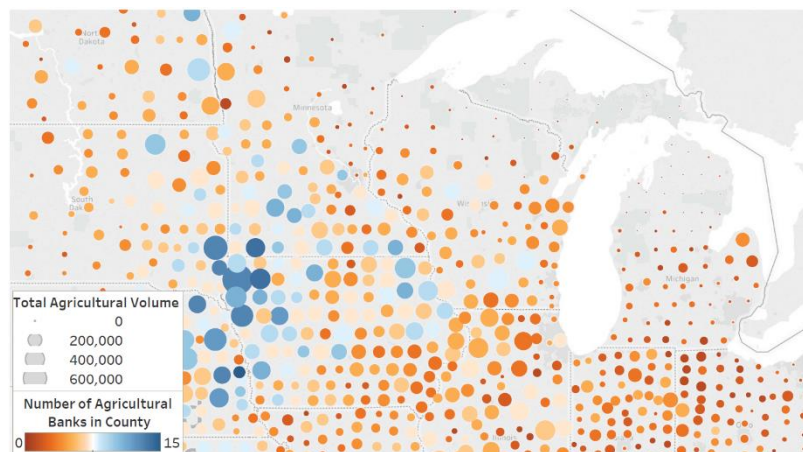
Number of Banks with > \$1Mn Loans in County



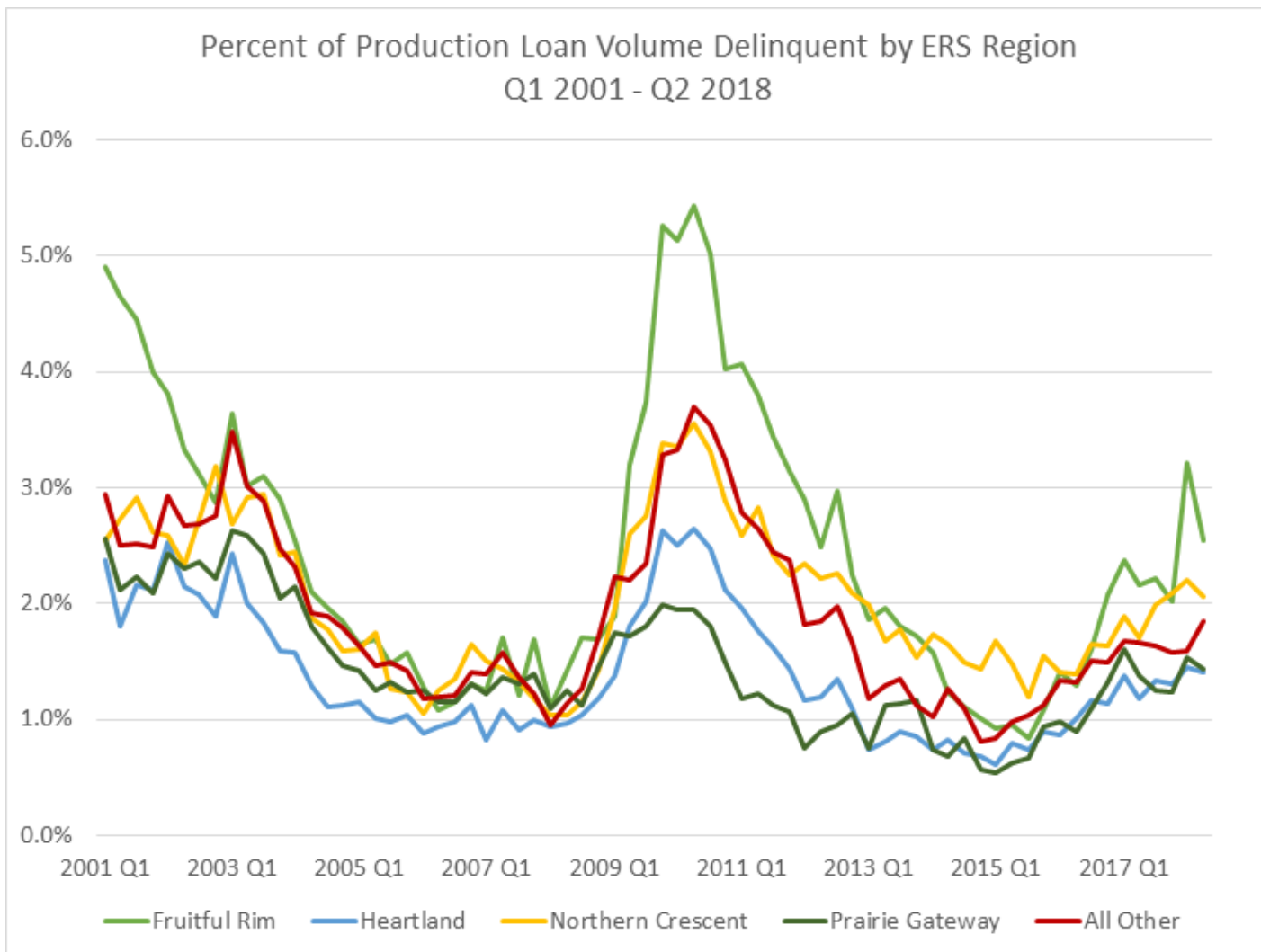
Can be used to check the robustness of agricultural credit markets across counties

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

Number of FDIC-designated Ag. Banks in County



Appendix: Production Loan Delinquencies by ERS Production Region



Appendix: Real Estate Loan Delinquencies by ERS Production Region

