

UNIVERSITY OF OKLAHOMA
GRADUATE COLLEGE

SOME NOVEL APPROACHES TO ECONOMIC PROBLEMS

A DISSERTATION
SUBMITTED TO THE GRADUATE FACULTY
in partial fulfillment of the requirements for the
Degree of
DOCTOR OF PHILOSOPHY

By

GABE LEBOVICH
Norman, Oklahoma
2022

SOME NOVEL APPROACHES TO ECONOMIC PROBLEMS

A DISSERTATION APPROVED FOR THE
DEPARTMENT OF ECONOMICS

BY THE COMMITTEE CONSISTING OF

Dr. James Hartigan, Chair

Dr. Tyler Ransom

Dr. Kenneth Haltman

Dr. Chunbei Wang

This work is dedicated to my parents Agi, Joseph, my brother Andy, my daughter Olive and Cheng. You have put up with my stubborn nature for all these years. This work reflects some of the fruit.

Acknowledgements

I would like to thank all my committee members. I am very grateful for Dr. Hartigan, who patiently and thoroughly analyzed all the economic problems I became interested in and helped me find my place in the larger world of economics. He is also the one who stressed the importance of empirical work to support my theories. I would also like to thank Dr. Ransom who valiantly accepted the task of advising me a mere six months ago. Without his superior intellect, skill, and patience the work laid out on the following pages would not exist. I would like to thank Dr. Chunbei Wang for her helpful comments throughout the years. I would also like to thank Dr. Haltman for his role as the outside member of the committee. Last but not least, I am thankful to Dr. Urquhart, Dr. Mott and Dr. Ho for their generous wisdom, patience and most of all friendship throughout the past two decades.

Contents

1	Understanding the role of social connections in the formation of political coalitions	1
1.1	Introduction	1
1.2	Literature Review	4
1.3	Data	7
1.4	Method	11
1.4.1	Level regressions	11
1.4.2	Shift regressions	13
1.4.3	Endogeneity challenge	14
1.5	Results	18
1.5.1	Level Regressions	18
1.5.2	Shift Regressions	19
1.5.3	2SLS Regressions	21
1.5.4	Heterogeneity Analysis	22
1.6	Conclusion	24
1.7	Appendix	26
2	Incorporating Bridging Social Capital into Social Capital Measures	36
2.1	Introduction	36
2.2	Descriptions and Literature	40
2.2.1	Definitions of Bonding and Bridging Social Capital	40
2.2.2	Current measurement of bridging and bonding social capital	42
2.2.3	Latest developments in network-derived social connectivity studies	44
2.3	Data	46
2.4	Method	48
2.4.1	Measuring Social Capital	48
2.4.2	Using Social Network Data to Measure Bridging Social Capital	50
2.4.3	Alpha centrality - the mechanics of combining the two types of social capital	52
2.4.4	Measuring the impact of social capital in various outcomes	53
2.5	Results	54
2.5.1	Descriptives and how alpha correlates with outcomes	54
2.5.2	Hate-group formation estimation using combined social capital	55
2.5.3	Social Capital Effects On Crime	57
2.5.4	Social Capital Effects On Other Outcomes	59

2.5.5	Centrality Gain From Increasing Network Radius	59
2.6	Conclusion	61
3	Heterogeneity in Child Quantity-Quality Trade-off: A Machine Learning Approach	69
3.1	Introduction	69
3.2	Data	73
3.3	Empirical Methodology	74
3.3.1	Heterogeneity in exposure to the One Child Policy	74
3.3.2	Instrumental Variable	75
3.3.3	Validity of IV	76
3.3.4	Local average treatment effect	78
3.3.5	Conditional local average treatment effect	79
3.4	Results	82
3.4.1	Local average treatment effect	82
3.4.2	Heterogeneity in local average treatment effects	82
3.5	Conclusion	88

List of Tables

1.1	Summary statistics	27
1.2	Level Regressions	28
1.3	Vote margin shift regressions. All connected counties.	29
1.4	Restricted network OLS and 2SLS shift regressions	30
1.5	Restricted network top and bottom quartiles shift OLS and 2SLS regressions	31
1.6	Split dependent variable along median 2SLS analysis	32
1.7	Split dependent variable 2SLS results for top and bottom quartile independent variables of interest	33
1.8	First stage of the 2SLS restricted network shift regressions	34
1.9	First stage results for 2SLS restricted network regressions with top and bottom quartiles	35
2.1	Eigenvector and alpha centrality status analysis of county level data of Oklahoma	63
2.2	Simple correlations between alpha centrality and outcomes for different values of attenuation parameter alpha	64
2.3	Effect of Social Capital on Hate Group Formation	65
2.4	Effect of Social Capital on Crime	66
2.5	Effect of Social Capital on Obesity and Unemployment	67
2.6	Radius Analysis by varying α	68
3.1	Heterogeneity in the OCP implementation	89
3.2	Summary statistics	90
3.3	Exposure to the strict OCP and sibling sex composition	93
3.4	2SLS: Height for age z-score and weight for age z-score	94
3.5	2SLS: relative education and school enrollment	95
3.6	Mean treatment effects by quartile	97
3.7	Heterogeneity in local average treatment effect by predicted treatment effect	98
3.8	Summary statistics by predicted treatment effect on HAZ	99
3.9	Summary statistics by predicted treatment effect on WAZ	100
3.10	Effect of number of siblings by mother's age at first birth	105

List of Figures

1.1	Intransitive Triads	26
2.1	Realtionship between geographically distant ties and weak links	62
3.1	Number of trees and median variance for boys	91
	(a) Boy HAZ median prediction variance	91
	(b) Boy WAZ median prediction variance	91
	(c) Boy relative education median prediction variance	91
	(d) Boy school enrollment median prediction variance	91
3.2	Number of trees and median variance for girls	92
	(a) Girl HAZ median prediction variance	92
	(b) Girl WAZ median prediction variance	92
	(c) Girl relative education median prediction variance	92
	(d) Girl school enrollment median prediction variance	92
3.3	Distribution of effects of number of siblings by gender	96
	(a) Distribution of effects of number of siblings on HAZ by gender	96
	(b) Distribution of effects of number of siblings on WAZ by gender	96
	(c) Distribution of effects of number of siblings on relative education by gender	96
	(d) Distribution of effects of number of siblings on school enrollment by gender	96
3.4	Health variable importance	101
3.5	distribution of effects of number of children on HAZ by parent’s education level	102
	(a) Distribution of effects of number of children on HAZ by father’s educa- tion level	102
	(b) Distribution of effects of number of children on HAZ by mother’s edu- cation level	102
3.6	distribution of effects of number of children on WAZ by parent’s education level	103
	(a) Distribution of effects of number of children on WAZ by father’s educa- tion level	103
	(b) Distribution of effects of number of children on WAZ by mother’s edu- cation level	103
3.7	distribution of effects of number of children on health by child’s age group .	104
	(a) Distribution of effects of number of children on HAZ by mother’s age at first birth	104
	(b) Distribution of effects of number of children on WAZ by mother’s age at first birth	104

Abstracts

In chapter 1, I investigate the role of social connections in the formation of political coalitions utilizing the county-to-county Social Connectedness Index (SCI) from Facebook. Using a generalized fixed effects model I estimate the impact of the SCI-weighted network vote margin on focal county vote margins. Democrat-leaning counties (above-median Democrat-net-Republican percentage) respond robustly to the influence of their social network in all years examined. Republican-leaning counties accelerated coalition formation in years after the 2012 election. Since voting for the US president happens at the same time in all counties, I also use a 2SLS analysis to deal with the endogeneity that stems from this simultaneous action. I construct the instrument for the endogenous network-wide vote margin variable by utilizing the network structure of the dataset. The empirical analysis in this paper confirms a high and growing level of political polarization in the US.

Chapter 2 creates a new measure of social capital that combines bridging and bonding social capital. New data from the Facebook user network makes this possible. I use the new measure to reassess the role of social capital on a variety of outcomes previously studied in the literature. I use the network topology of social connections and within-county associational membership levels to measure county-level social capital. This measure allows one to combine two different types of measures into a comprehensive one. Since bonding and bridging social capital are by nature different, the ability to measure these two values from different sources and then combine them improves my measure of social capital over previous ones. Results show that my measure is a superior predictor of a variety of social and economic outcomes.

Chapter 3 is coauthored with Cheng Ma. In chapter 3, we test the effect of additional children in a family on health and educational outcomes. Estimating this effect is complicated by the endogeneity of family size. We use the variation in the severity of the effect of the one-child-policy in China to extract exogenous variation from the China Health and Nutrition Survey in the same country. After finding a negative effect of family size on health and educational outcomes of children we use a newly developed machine learning approach. Generalized random forests allows us to look at the heterogeneity in treatment effects in the quantity-quality trade-off. We find robust negative treatment effects of additional children on health outcomes but only mild effects on educational attainment. The machine learning algorithm finds mother's age and parent's education level play a large role in the negative quantity-quality trade-off. Pinpointing the factors that exacerbate the negative effect of additional children on child quality can aid future policy decisions.

Chapter 1

Understanding the role of social connections in the formation of political coalitions

1.1 Introduction

Traditionally researchers have looked at the impact of economic and demographic factors on voter preference. However, voters themselves may exert influence on each other leading up to an election. The recent availability of social network data from websites such as Facebook and Twitter allows the examination of the impact of comprehensive social networks on voting decisions. What is the impact of social networks on voting decisions in presidential elections? Is there partisan heterogeneity in the impact of social networks?

This study finds that distant networks influence local voting using the Social Connectedness Index provided by Facebook Corporation. For a one percentage point increase in the average network county presidential election vote margin there is a 1.04 percentage point increase in the local county's vote margin.

To find the impact of comprehensive, country-wide networks on focal county vote decisions in presidential elections I combine county-to-county aggregate Facebook connections with county-level electoral outcomes. The Social Connectedness Index from Facebook provides the intensity of connections every county in the US has with its social network. The level of importance that each friend county's vote margin provides to the comprehensive influence of the social network is measured by the SCI. I measure the effect of network-wide social network influence on voting using an OLS model. I control for county-level social, demographic and economic factors in the model. In addition to the weighted network-wide vote margin I also look at the impact of the weighted four-year vote margin shift on focal county vote margin.

The latter research design from above allows a partition of the network into partisan halves. The impact of Democrat-leaning social network counties is larger and more consistent. The effect of the Republican portion of the social network on Republican-leaning focal counties is to increase the vote margin in favor of the Republican candidate. However, the effect of the Republican voting social network on Democrat-leaning focal counties is to increase net votes in favor of the Democratic candidate. The Democrat part of the social network increases vote margins in favor of the Democratic candidate whether the focal county is Democrat-leaning or Republican-leaning.

There have been several studies so far using Facebook data to look at the impact of social networks on decision making.¹ This study adds another chapter in the recent utilization of Internet social network site data to study the impact of comprehensive social networks on local decisions. It specifically examines the impact of social networks on local voting decisions at a county level while controlling for social, demographic, and economic characteristics. The research design is unique because it can pinpoint network effects despite there being no time variation in the activity of interest: elections. The simultaneity of actions between the

¹Bailey et al. (2018a) look at the effect of distant networks on housing investment decisions, Wilson (2019) examines how social links influence EITC claiming, Bailey et al. (2018b) look at the relationship between social connectedness and a number of economic factors.

network and the unit of analysis, the county, create the potential for endogeneity as explained by Manski (1993). I use the network structure of the dataset to construct instrumental variables to overcome this problem. Using county-to-county social connections provided by Facebook this is the first study to examine the effect of the comprehensive, countrywide network on local county voting patterns.

There are two main sources of variation in the empirical analysis. A county X is connected to different counties than county Y . Also, even when X and Y have connections to the same county the level of connections will be different. County X might have more friendship links to county Z than county Y does. These two sources of variation generate variation in the social networks each county faces. The level of connections are used as weights to find the weighted average of network-wide vote behavior that each focal county is exposed to. Then I look at the effect of this variation on focal county vote margins.

I find robust local voting outcome responses to social network influences in the four most recent presidential elections. Whether the network influence is quantified by vote margin changes (shifts) from one election to the next or network vote margins there is a statistically significant relationship between casting presidential votes in a focal county and the voting in its social network. In most elections the relationship is positive, with network-wide moves into certain political directions translating to similar movements locally as quantified by casting votes for Republican vs. Democrat presidential candidates. In the specification where the network vote margin is used I find a positive relationship on local voting outcomes in presidential elections in a panel specification as well as the 2008, 2012, 2016 and 2020 elections individually. The election-to-election network-wide vote margin change is positive and significant in all years except 2012. In 2012 the effect of the Republican portion of the social network was a relative increase of local votes to the Democratic presidential candidate. A 2SLS model implemented to deal with endogeneity confirms my findings in the OLS models. The paper also examines heterogeneity over Republican vs. Democrat focal counties and also looks at treatment intensity (i.e. network connectedness intensity).

These findings reveal the effect of social network influence in presidential elections. They also reveal the effect of networks on decision making in general. This effect is gleaned from Internet-based social connections between people. However, since in most cases these connections reveal already-established connections between people, studies using data from social network ties provide a window into social network influence as it relates to real-life interactions beyond the Internet.

The paper is organized as follows. The next section will examine the literature on utilization of social networks to study economic and social questions. Section 3 explains the data used in this study. Section 4 looks at the method used to analyze the dataset. I discuss the results in section 5. Section 6 concludes.

1.2 Literature Review

This study supplements traditional determinants of voting outcomes such as demographic, social and economic factors², by adding social network influences³ to the analysis. Therefore, in addition to examining the literature on voting outcomes it is helpful to look at the emerging literature of social networks. The recent availability of social connectivity graphs has been made possible by data collected by social network websites such as Twitter and Facebook, among others.⁴

There is no shortage of studies looking into a wide range of determinants of presidential elections outcomes.⁵ Other studies detail the role of the Internet in recent electoral out-

²For a detailed description of each determinant see Adkisson and Peach (2018)

³Bailey et al. (2018b), Wilson (2019) detail social network influences in their studies.

⁴Putnam et al. (2000) claims that Internet connections are unable to foster the formation of links that build trust and reciprocity; however he was writing this before the time of social tie-building websites such as Facebook and Twitter. Sajuria et al. (2015) examine the social capital building potential of social network websites and conclude that there is evidence of significant build-up of bonding social capital, such as the links between close friends and family, and evidence of conditional bridging social capital, such as weak ties that occur between acquaintances.

⁵Enke (2020) analyzes the economics of moral decision-making presidential elections, Adkisson and Peach (1999) look at county-level border region heterogeneity, Flaxman (2018) looks at the county-level effect of both economics and white identity in the 2016 election

comes.⁶ Not only has the Internet become more important at the expense of traditional print media and TV, but Internet connections between people provide researchers with promising datasets to study the effect of social networks on individual decisions.

The above studies use several common controls that their authors deem relevant determinants of presidential election outcomes. These determinants can be grouped into social, demographic and economic factors.⁷ The social factors include percentage of people of a certain race, ethnicity or religion in addition to educational outcomes, such as the percentage of population with a certain educational attainment. Population density and percentage of the population of a certain age are typical demographic determinants. The economic determinants of voting outcome are proxied by per capita income and unemployment levels as in Flaxman (2018). The county-level variation of these factors plays an important role in choosing a president.

I use the social connectivity data from Facebook to build a county-to-county social connection variable. However, before I discuss the papers that also use Facebook data it is helpful to mention two earlier studies that use social network data and look at the relationship between network structure and socioeconomic characteristics. Sajuria et al. (2015) look at the dynamics of social capital formation in Internet-aided political movements such as Occupy Wall street, the Chilean presidential election, and the UK IF (end hunger) campaign. The authors look at Twitter data from which they glean social connectivity that forms via the Internet at these events. As mentioned above they find evidence of social capital formation. There is evidence of bonding social capital, such as links between close friends and family members, when they look at clustering of social ties. They use the Burt (1992) theory of structural holes to test for brokerage which is another way of calling bridging capital (social ties between acquaintances). By looking at how much of an individuals' connections are concentrated in a small group they are able to determine how much ability

⁶Allcott and Gentzkow (2017) look at the role of fake news in the 2016 election cycle, Allcott and Gentzkow (2017) look at the effect of Internet use in general in the 2016 elections

⁷The variables in these three groups of vote determinants are discussed by Adkisson and Peach (2018).

or access these individuals have to bridging structural holes, i.e. how easy it would be for them to form links to other connected groups. They find bonding social capital is organically created in all three movements but only the movements with professional brokers are able to reduce network constraint and create bridges between the different groups. Their study is one of the first to examine real world comprehensive social network structures.

Eagle et al. (2010), using cell phone data from the United Kingdom, look at the relationship between social diversity and socioeconomic status. Their socioeconomic variable is the Index of Multiple Deprivation (IMD) composed of income, employment, education, health, crime, housing, and the environmental quality metrics. The social and spatial diversity variables use cellphone data. They find a positive correlation between the IMD and both social and spatial diversity. In addition, they determined there is a negative correlation between time spent communicating on cell phones and IMD. The authors note that they could not establish causation between these variables because of a lack of longitudinal data.

Recently, more and more social network sites have been providing valuable data about real world social ties and networks. The social network dataset used in the present paper comes from Facebook. Bailey et al. (2018b) were the first ones to analyze the dataset. In addition to basic Facebook connectivity statistics their paper looks at friendship link elasticities with respect to distance. In a subsequent paper the same authors find significant effects of house price changes in a person's social network on their homeownership decisions using individual level connection data (Bailey et al. (2018a)).

The Facebook dataset also reveals correlations between trade flows and SCI. While this mechanism is in accordance with established trade theories the SCI measure is able to explain more of the variation in trade flows than distance. Social connectivity between counties is able to overcome informational and cultural frictions. Bailey et al. (2018b) also find correlation between SCI and innovation by looking at patent citations. In addition, they find a relationship between connectedness and migration which they hypothesize is caused by the higher willingness to move to a county where one has friends or family who live there.

There have been a handful of other studies that use SCI from Facebook and other social connectivity data. Wilson (2019) uses the SCI dataset to look at Earned Income Tax Credit (EITC) take-up. He finds that while out of state EITC adopters don't have influence on EITC take-up in the focal state they do influence the level of business income in tax forms. People adjust their non-wage income to claim maximum EITC benefits as a result of out-of-state social network influence as measured by the Facebook friendship links.

Consistent with Bailey, Wilson (2021) finds that the level of county-to-county migration drops significantly when the two counties in question are separated by state borders. Since friendship links also drop off at state borders, the reduction in migration might be the work of information frictions, personal connections or different state identities. This reduction in migration levels at state borders hinders the economic resilience of counties sharing a state border in the face of economic downturns such as the Great Recession.

1.3 Data

To estimate the impact of social networks on county-level election outcomes I use the Social Connectedness dataset from Facebook. The dataset aggregates Facebook connections to the county-level. I combine it with county-level voting results in the 2008 - 2020 presidential elections provided by the MIT Election Lab. I am also able to obtain, from the American Community Survey, county level aggregations for all the control variables: unemployment rate, per capita income, population, percent of population over 65 years of age, percent of population Hispanic or Black, and percent of population with a bachelor's degree who are older than 25 years of age. I use the Facebook social connectivity measure as weights to add up network-wide voting outcomes. Variations on this method constitute my principal variable of interest. The response variable is vote margin in the 3105 focal counties (including 130 county equivalents). These 3105 counties are the unit of analysis in this study. Before describing in detail the datasets I use in the analysis it is useful to look into the novel dataset

from Facebook.

According to Bailey et al. (2018b), there are 239 million Facebook users in the US, highly representative of the US population.⁸ The establishment of a friendship connection in Facebook requires the consent of both parties. There is a limit of 5000 links any one user can establish, limiting the number of popularity-signaling “vanity links”. The connections between Facebook users have been aggregated to the county level for all of the 3105 counties. While Facebook only released the county-pair connectivity values for 2016 they claim that usage of the platform is fairly stable overtime.

The Social Connectivity Index measures the number of Facebook connections between all county pairs in the US.⁹ The variation of normalized county pair connections is identical to the variation in the real number of connections between all county pairs. Facebook also released the Relative Probability of Friendship measure. This is the number of friendship links between county-pairs divided by the product of total Facebook account holders in the two counties. This measure reflects the number of friendship links established out of all possible links that may be formed between any two counties. In this manner the measure shows the probability of friendship connections between the two counties and is not distorted by the relative populations of the two counties. While the present study does not utilize the Relative Probability of Friendship the measure provides a useful insight into the nature of the Facebook connectivity dataset.

Bailey et al. (2018b) also report summary statistics of the connectivity in the SCI dataset. In some of the analysis that follows I restrict the network to out-of-state counties only. For this reason it helps to look at already-established relationships of friendship links with distance. The intensity of friendship links is strongly dependent on distance. The elasticity of friendship links with respect to distance is -2 within a 100 miles and is attenuated to -1.2

⁸Facebook only counted accounts that have been logged into within the last 30 days of the snapshot of their connections release.

⁹For data security purposes this measure is normalized to have the largest value of 1000000. The county with the largest Social Connectivity Index (SCI) that is normalized to 1000000 is the number of connections that Los Angeles county has with itself. In other words, this measure is the normalized one of all the real connections that Losangelinos have with other Losangelinos.

for distances longer than 200 miles. On average, 62.8% of friendship links occur within a 100 miles. However, there is large variation in this value from 46% at the 5th percentile to 76.9% at the 95th percentile. This variation in the number of distant friendships is closely correlated to a number of socio-economic variables. Counties with more distant ties tend to be better educated, have higher incomes and have higher social capital.

Facebook does not provide individual-to-individual connections in the dataset it shares with the public. The SCI dataset includes county-to-county scaled Facebook connections. The number of Facebook connections that all individuals who live in county c have with all individuals who live in county j , for all counties j . No additional data is provided by Facebook. Therefore the influence measures constructed above measure general propensities between the connected county populations. The Democrat exposure measure (aggregate influence on county c by all counties j where the margin of Democrat minus Republican vote margin went up over four years) includes connections between Democrats in county c and Democrats in county j and Democrats in county c with Republicans in county j , in addition to the other two permutations. However, overall the change in Democrat exposure, which I denote $friend_d_exp$, increased in county j over four years. Similarly for $friend_r_exp$ indicates all the SCI weighted decreases to counties that county c is connected to where the margin of Democrat net Republican decreased in four years. To delineate a clearer picture of how voting changes in friend counties might influence a focal county I also look at the influence of only those friendship counties where the change in the vote margin of Democrat over Republican over four years is in the top quartile (for Democrat influence $friend_dd_exp$) or in the bottom quartile ($friend_rr_exp$).

The rest of the independent variables are county level measures that influence how people vote as described in the literature review above. These variables will be used as controls in the model. All of the control variables come from different waves of the American Community Survey. The ACS is an ongoing survey that tracks county-level (and other geographic level) data on employment, educational attainment, and demographics annually in the US. For 2008

I use the 2009 ACS 5 Year Data to calculate these measures. The values for these variables are averaged over 5 years from 2005 to 2009. The 2012 wave of the control variables are averages from 2008 through 2012 also from the ACS. For the 2016 controls the variables are averages from 2012 to 2016 and for 2020 the values for years 2015 through 2019 are averaged. The vote data was obtained from MIT Election Lab.

Table 1.1 contains the means and standard deviations of the variables used. This table lays out the ranges of the relevant variables. It aids in visualizing the changes in the variables over the four elections considered. The Democrat net Republican vote margin goes from -15% in 2008 to -31% in 2016 and remains there in 2020. The mean of the variable that measures the vote margin of the network, *friend_exp* is -0.078 across all years. It is 0.008 in 2008 and -0.106 in 2020. This means the SCI-weighted vote margin that the average focal county faced through its' social network was around zero in 2008 and became negative in 2020. The mean of *friend_d_exp* shows the SCI-weighted vote margin change faced by that part of a county's network where the vote margin change was positive from one election to the next. The mean of this variable is 0.02 across all years. The mean of *friend_r_exp* across all years is -0.04. The difference in means in these two measures is consistent with the vote margin declining over the elections considered.

The mean across focal counties of the "share of population of 65 year olds" variable increases from around 16% in 2008 to 18% indicating an aging population. The mean value across all counties of the share of Blacks and Hispanics is fairly steady over all election years considered. The mean of "percent of population with a bachelor degree who are 25 years of age or older" goes from 18% in 2008 to 27% in 2012. It goes down to around 20% for the remaining two election years. The mean unemployment rate is between 5% and 8%. The mean of incomes across all counties rises from 22000 to 28000 from 2008 to 2020.

1.4 Method

With the data in hand, I now model how social connections might influence election outcomes in addition to demographic and economic variables previously analyzed in the literature. I first illustrate this influence by OLS regression. I then make use of an alternative panel data model as well as instrumental variables.

1.4.1 Level regressions

I first use OLS regression to examine the effect of social networks on county-level vote margins in presidential elections. To find the effect of social networks on voting outcomes I add up vote margins in counties the focal county has Facebook friendship ties with. I use the strength of ties as weights for the friendship county's vote margin. In this manner the vote margins in those counties that the focal county has relatively more connections to play a more significant role in the total network vote margin summation. A natural starting point in the regression analysis is an OLS model.¹⁰

$$Vote\ margin_{cst} = \beta friend_exp_{ct} + \gamma X'_{ct} + \theta_s + \delta_t + \epsilon_{cst} \quad (1.1)$$

where c indexes counties, s indexes states, and t indexes years.

The dependent variable is the percentage of votes in a county that went to Democrats minus the percentage that went to Republicans. That is,

$$Vote\ margin = \%vote\ to\ Democrat - \%vote\ to\ Republican \quad (1.2)$$

The independent variable of interest is $friend_exp_{ct}$ which is the margin of the percentage of votes to the Democratic candidate minus the percentage gathered by the Republican one in each of the above elections in friendship counties times the connection strength or number of friendship links. X represents demographic and economic variables. Since it is useful to

¹⁰This strategy has been implemented by Wilson (2019) in estimating the effects of social networks in EITC claiming.

examine coalition formation over time, initially I separately examine the elections of 2008, 2012, 2016 and 2020. In the panel regression, time fixed effects are indicated by δ_t , while state fixed effects are θ_s .

The links that county c has with its friendship counties measure the intensity of social connection between the two counties. To find the aggregate influence that all counties, connected to c exert on it, the margin of votes as described above is multiplied by the number of connections that county c maintains with county j , and then all these values for every (c, j) pair are added up. The SCI can be thought of as a weight used to add up all vote differences in all the friendship counties. The SCI weights are calculated as the proportion of links to the particular friend county out of the total out-of-county links of the focal county.

$$friend_exp_{ct} = \sum_{j=1}^{N-1} \left(\frac{SCI_{c,j}}{\sum_{k=1}^{N-1} SCI_{c,k}} \times Vote\ margin_{jt} \right) \quad (1.3)$$

The rest of the independent variables found in vector X are operationalizing other factors that influence electoral outcomes. Economic controls such as the log of income and the unemployment rate account for between-county variation in economic conditions. Demographic variables such as the percentage of population 65 years or older focus on the variation between counties in the number of people who have seen the US in more high-growth times after World War II. In addition, this demographic tends to have a more closed network. According to Bailey et al. (2018b) the clustering coefficient for this group is 12.5% (versus 9.4% for the 25 to 35 year olds) meaning 12.5% of people in this group's friends are friends among themselves. The coefficient on this variable is expected to be negative reducing the Democrat minus Republican percentage in a focal county.

The logarithm of the population accounts for rural/urban areas and for population declines due to economic instability. Hence the coefficient on this variable is expected to be positive. Social influence variations between counties are operationalized by the percentage of black voters, percentage of Hispanic voters and the percentage of bachelor degree holders over 25 years of age in a county. Blacks and Hispanics typically vote Democrat; however

the data used is aggregated to the county level. While black and Hispanic voters may vote for the Democratic candidate, their neighbors might not — especially if they are concerned about racial diversity in their locales. State fixed effects are also included in all models since states have different cultures and some have different electoral laws. Furthermore Bailey et al. (2018b) found strong drops in the number of Facebook friendships occurring at state borders possibly due to shared state identities, economic opportunities, shared legal systems or institutions.

1.4.2 *Shift regressions*

The level regressions don't include county fixed effects. First-differencing the independent variable over the four-year election-to-election period gets rid of some county-level unobserved heterogeneity. In addition, I want to measure the effect of network-wide vote changes over time on focal county vote margins. This strategy has been proved useful in a social network setting by Bailey et al. (2018a) in measuring the effect of network-wide house price changes over time on own house purchase decisions. Bailey et al. (2018a) refer to the shift in network-wide house prices as the house price experience over the time period considered that the focal person is exposed to over his network. The shift in vote margins from one election to the next also allows to bifurcate the independent variable of interest into influence by counties where the shift over four years was in the advantage of Democrats and influence by counties where the shift in vote margins was in the direction of Republicans. In counties where the four-year shift in vote margin ($Democrat\% - Republican\%$) was positive is the Democrat influence: $friend_d_exp_{ct}$. Conversely, where this shift over four years is negative it is the Republican influence: $friend_r_exp_{ct}$. In this specification the number of elections drops to three: 2012, 2016 and 2020. I use the nomenclature from Flaxman (2018) and call this a vote-switch model.

$$friend_d_exp_{ct} = \sum_{j=1}^{N-1} \left(\frac{SCI_{c,j}}{\sum_{k=1}^{N-1} SCI_{c,k}} \times (Vote\ margin_{jt} - Vote\ margin_{j,t-4}) \times \mathbb{I}\{Vote\ margin_{jt} - Vote\ margin_{j,t-4} > 0\} \right) \quad (1.4)$$

$$friend_r_exp_{ct} = \sum_{j=1}^{N-1} \left(\frac{SCI_{c,j}}{\sum_{k=1}^{N-1} SCI_{c,k}} \times (Vote\ margin_{jt} - Vote\ margin_{j,t-4}) \times \mathbb{I}\{Vote\ margin_{jt} - Vote\ margin_{j,t-4} < 0\} \right) \quad (1.5)$$

After partitioning the independent variable of interest the new model is this:

$$c_marg_{ct} = \beta_1 friend_d_exp_{ct} + \beta_2 friend_r_exp_{ct} + \gamma X'_{ct} + \theta_s + \delta_t + \epsilon_{ct} \quad (1.6)$$

c_marg_{ct} is the margin of county-level votes in year t as in the level regression. $friend_r_exp_{ct}$ is negative for all years since it is the interaction of the positive SCI values and negative vote margin changes. Similarly, $friend_d_exp_{ct}$ is positive for all election years and the pooled model as well. For additional summary statistics for $friend_r_exp_{ct}$ and $friend_d_exp_{ct}$ please consult table 1.1.

1.4.3 Endogeneity challenge

Presidential elections happen on the same day for the whole country. There is no time variation in the event which makes isolating a channel of causality difficult. As voters in one county get influenced by voters in the counties they have social connections with, they in turn themselves exert influence on these voters in these counties. This simultaneity of causality has been termed the reflection problem by Charles Manski. He explains the reasoning for this nomenclature on the example of an alien who, seeing a person staring at his reflection in the mirror, is unable to discern if it is the reflection making the person move or the other way around. The reflection problem in general arises when a researcher tries to determine the influence of the mean behavior of a group on an individual's behavior who is part of this

group.

Manski's Identification Conditions

There are two social effects laid out by Manski (1993). One is endogenous effects — the propensity of an individual to be influenced by the average behavior of her group. The other is exogenous effects or the propensity of an individual to be influenced by the average exogenous characteristics of his group. In addition to these foregoing social effects there is a correlation effect — the propensity for an individual to behave similarly to her group because they are embedded in the same institutional environment and/or have similar individual characteristics. Upon deriving the reduced form equation Manski (1993) observes that the endogenous effect is not identified if:

1. The variable that defines the groupings (x) is a function of the exogenous variables (z).
2. $E(\text{group variable} \mid \text{exogenous variable})$ does not vary with the exogenous variable.
3. $E(\text{group variable} \mid \text{exogenous variable})$ is a linear function of the exogenous variable; the grouping variable is a linear function of the exogenous variable.

It is necessary to examine the dataset to see if its structure allows identification. The providers of the dataset note that the connections between friends in the Facebook dataset reflect real life friendships. Some of these social connections are the result of historical migratory patterns. Therefore there is little concern that we are dealing with endogenous network formation at the county level. Among the examples that Bailey et al. (2018b) list are the connections between Oklahoma and Kern County, California. This connection arose from the migration of Oklahomans to California during the Dust Bowl in the 1920's. Another example they list is the migration of freed slaves from Mississippi to the Chicago area. The connections between these two pairs of places persist into present times. If most connections in the dataset were the result of recent resettlements, the network would be itself endogenous. The connections between counties would be in large part the result of the

relevant people's characteristics. Since this is not a concern the only remaining endogeneity that must be dealt with is simultaneity: the reflection problem described above.

Intransitive Triads

Bramoullé et al. (2009) refine Manski's model for networks. This extension looks at the identification of endogenous effects in a network setting. The presence of intransitive triads (a simple intransitive triad is where A is friends with B and B is friends with C but C is not friends with A) in general allows for the opportunity for identification.

Making a connection with panel data models, the authors observe that instead of time lags of exogenous variables as used in panel data, the exogenous variable's of neighbors of neighbors (spatial lags) can be used as instrumental variables for identification of endogenous social effects with network data. In general if the focal node is excluded and group size varies, identification is obtained. In cases where the focal node is included in the group, identification fails. For example, in a classroom setting if a student is included in the estimation of social effects of her group, her mean family background of her friends is equal to the mean family background of her friends' friends. In the present paper the focal county is excluded from the group. Intransitive triads are present. Group size varies. This data structure allows the use of instrumental variables to deal with the endogeneity problem.

Instrumental Variables Solution to Endogeneity

Sunder et al. (2019) use the method developed by Bramoullé et al. (2009) to develop an instrument using intransitive triads to estimate the effect online ratings on herding behavior.

Similarly to the model developed by Bramoullé et al. (2009) and implemented by Sunder et al. (2019) I employ an instrument derived from intransitive triads. The key to the development of the instrumental variable is the presence of counties that have links to the focal county's friend county but do not themselves have any significant ties to the focal county.

So if focal county A has links to counties B , I will utilize all counties C that are connected

to counties B but not to county A . (Please see figure 1.1 for a graphical explanation.) Then I use the mean of the exogenous variables of counties C to instrument for the endogenous variables $friend_d_exp_{ct}$ and $friend_r_exp_{ct}$. The first and second stages are as follows:

$$friend_d_exp_{cst} = \alpha_1 Z_{ct} + \gamma X'_{ct} + \theta_s^1 + \delta_t^1 + \epsilon_{cst}^1 \quad (1.7)$$

$$friend_r_exp_{cst} = \alpha_2 Z_{ct} + \gamma X'_{ct} + \theta_s^2 + \delta_t^2 + \epsilon_{cst}^2 \quad (1.8)$$

$$c_marg = \beta_1 \widehat{friend_d_exp}_{cst} + \beta_2 \widehat{friend_r_exp}_{cst} + \gamma X'_{ct} + \theta_s + \delta_t + \epsilon_{cst} \quad (1.9)$$

In this specification the dependent variable remains the difference in the percentage of votes to Democratic presidential candidate and the percentage of votes to the Republican presidential candidate in a focal county c (c_marg). Also there are two independent variables of interest. As above, one is composed of the change in the margin of percentage of votes for Democrat minus percentage of votes for Republican presidential candidate from one election to the next if the change is positive. These are the counties where over four years the margin of Democrat net Republican percentage of votes moved positive — more people voted Democrat in year $t + 4$ than in t . β_1 and β_2 would estimate the Local Average Treatment Effect if the instrument was binary. However in the present case the instrument is a continuous variable. β_1 and β_2 represent the continuous analog to LATE, however they do not estimate LATE.¹¹

The statistics of the validity of the instrument as revealed by the first stage are published in table 1.8. The F-statistics exceed the standard threshold, indications a strong instrument. A Wald test and a Conditional Likelihood Ratio test confirm this. Please see appendix 1.8 for details. The instrument meets a strict exclusion restriction as discussed in section 4.3.2 .

¹¹Blandhol et al. (2022) note that without a rich covariate set and monotonicity correct first stage the coefficients from 2SLS cannot be given a weakly causal interpretation.

1.5 Results

In this section I present the estimates of the level and vote shift regressions. First I discuss the level regressions. Here I look at the results pooled across all presidential elections from 2008 through 2020, then for each election year individually. I look at the effects of the comprehensive network and out-of-state network. Subsequently, in the same fashion, I look at the panel estimation for vote shift regressions pooled and one by one. To address endogeneity concerns, I also estimate a 2SLS equation using the instruments developed above. Next I examine treatment intensity by looking at the top and bottom quartile of the independent variable. Finally, I look at the network influence separately on focal counties where the vote margin is above the median, Democrat; and below the median, Republican. The coefficients from these regressions confirm my research hypothesis. There are consistent effects of the social network on voting outcomes. Furthermore, there is heterogeneity in social network influence on Republican and Democrat counties.

1.5.1 Level Regressions

The first specification, estimating equation (1.1), looks at the dependence of vote margin in a county on the average friendship county vote margin, weighted by the SCI index. There are 4 election cycles included in the first specification: 2008, 2012, 2016 and 2020. Because of the inherent panel structure of the data in the first specification, all years are included in the regression. I include year and state fixed effects. Year fixed effects control for unobserved heterogeneity in part stemming from the various candidate pair-ups in the different elections. State fixed effects control for correlated effects such as culture and legal system.

Evident from column 1 of 1.2, the coefficient on network exposure is positive and significant at the 1% level. A unit increase in friend exposure results in a 1.04% increase in the percentage of county-wide vote to Democratic candidate minus the percentage county-wide vote to the Republican candidate. The coefficient is the same sign and of similar magnitude

as the one on “percent of population over 25 with a bachelor degree.”¹²

The within-state county connections give rise to inherently endogenous networks since people in a state share a unique legal system, culture and possibly climate.¹³ For this reason I look at exposure for the out-of-state social network in the next specification. This variable is constructed similarly to friend exposure; however, only out-of-state county votes are added up, weighted by the SCI friendship weights. The coefficient of interest increases to 1.39%, as seen in in column 6, and is significant at the 1% level. The increase in friend county influence when only out-of-state counties are included is consistent with the Eagle et al. (2010) paper mentioned above.

The above specification hides election year heterogeneity that is only available in individual election year analysis. In the next step I look at the effect of the social network on local voting decisions for each election year individually. State fixed effects remain in the equation to vouch for state-level time invariant characteristics such as culture and legal system. Columns 2 through 5 show the coefficients are 1.5% 1.31% 0.98% and 0.95% for 2008, 2012, 2016, and 2020, respectively. From the above specifications it appears that there was more coalition building in 2008, however it continues through the 2020 election. A cursory look at the decreasing coefficients unveils that the scope for future coalition building is limited. The annual coefficients for out-of-state exposure are 1.5% 1.48% 1.09% and 1.07% for 2008, 2012, 2016 and 2020 respectively. The rest of the analysis looks at network influence separately for the Republican part and the Democrat part of a county’s social network.

1.5.2 *Shift Regressions*

The shift regressions in equation (1.6), look at the local effect of the four-year change in the vote margins in counties the focal county has social connections to. To accomplish this, the independent variable of interest in one election year is subtracted from this variable in the next election year. This exercise cuts down the dataset to 3 election years — 2012, 2016

¹²The coefficients on the control variables are consistent with the literature.

¹³Facebook reports sharp drops in the elasticity of friendship links to distance at state borders.

and 2020. The SCI-weighted friend-county vote margin shift allows me to split the data set into counties where the shift was favorable to the Democratic party (*friend_d_exposure*) and to the counties where the shift favored Republicans (*friend_r_exposure*). The value of Democrat friend exposure is comprised by vote shift values in friendship counties where the election-to-election vote shift in Democrat margin was above the median value, weighted by the SCI. Similarly for the Republican friend exposure, only those counties' vote margin shift values are added up where this value is below the median. The coefficients in table 1.3 reveal interesting information about the role of coalition formation in recent elections. For the panel of three elections, as the Democrat friend exposure variable increases one unit vote margin increases by 1.98 percentage points for Democrats. Republicans cut into the Democrat net Republican difference 0.4 percentage points across all three years.

The panel data hides significant election year differences since there were different candidate pair-ups in the 2012, 2016, and 2020 elections. Columns 2 through 4 look at network effects in individual elections. Barack Obama's campaign relied heavily on social media and the Internet in general to get people to the polls. The influence of Democrat counties on a focal county on average was a 7.4% increase in local vote margins in 2012. In 2016 Hillary Clinton's messaging failed to motivate a similar social influence campaign with only 0.53% increase in Democrat over Republican margin for a one unit increase in friend county exposure to counties where four year shift in Democrat minus Republican vote margin was above the median. Since the independent variable measures 4-year shifts in vote margins and is split into two variables, one where the shift is positive and one where it's negative, the coefficients can be used to measure political polarization. The largest level of Democratic polarization occurred in 2012, while Republican polarization is more modest but growing over the elections considered.

It is puzzling that there is the positive and significant 0.84% (significant at 0.01%) coefficient for the Republican friend county exposure in 2012. This coefficient means a unit increase friend county exposure to counties where from 2008 there has been a below median

shift in Democrat minus Republican (an increase in votes to Republicans) increased the margin of Democrat over Republican by 0.84%. In 2016 the coefficient for Republican social network exposure is -1.15%, indicating that the the Trump campaign successfully capitalized on the power of the Internet to build a coalition to edge out Hillary Clinton to claim a plurality in the 2016 election. In 2020 the coefficient on Republican social network exposure is even higher in magnitude at -3.4%. For the 2020 election, Barack Obama’s former running mate Joe Biden using similar campaigning tactics as the Democrat campaign in 2012 and was more successful in rallying voters. The coefficient on Democrat network exposure increased to 4.05% — not as high as Barack Obama’s but about 8 times higher than Hillary Clinton’s.

1.5.3 2SLS Regressions

As discussed above there is a possibility of endogeneity in the above analysis. To overcome it the endogenous variables *friend_r_exp* and *friend_d_exp* are instrumented (indicated in equation 1.7) by the average exogenous characteristics of the friend counties of friend counties who do not have any significant connection to the focal county. Since it is rare for a county to not have any connections to any other county, the restriction of only connections to counties with above median number of Facebook friends are counted. If a county has lower connections than the median these connections are not counted. This is only a minor restriction since the largest number of connections are skewed toward nearby counties. These connections make up the largest number of links as do within-state connections. Most counties have connections to some counties out of their state but only to a few of them in any substantial way. Therefore, for the 2SLS specification, only the links in the county’s top 50 percentile are counted. This restriction only gets rid of connections with very few links and clears up the presence of intransitive triads.

First in table 1.4 are OLS results for the restricted dataset. These coefficients are not substantially different from the previous specification confirming that the effect of restricting

the dataset to the upper 50th percentile of connections will not alter any results. The coefficients from the 2SLS analysis ¹⁴ reveal the endogeneity biased the estimates downward. The magnitude of the 2SLS coefficients is larger while retaining some proportionality to the OLS coefficients. As mentioned above the coefficients in the 2SLS model estimate a continuous analog to LATE. Always takers are mixed in the estimation with compliers. Since always takers are negatively weighted in the estimation the coefficients will not estimate LATE.

The coefficient on *friend_d_exp* is 219.1 for 2012 which is about 30 times larger than the OLS estimated coefficient. Interestingly, the 2SLS coefficient for Clinton is larger at 37.14% than Biden’s coefficient at 22.76%. It is unclear why this is the case for the 2SLS model but not for OLS. The coefficients for the Republican friend exposure are larger in magnitude but proportional to the OLS coefficients with one notable exception. The positive coefficient for 2012 (increase in Democrat over Republican margin) loses significance in the 2SLS analysis.

1.5.4 Heterogeneity Analysis

To see if the above results are driven by counties that have shifted the most towards choosing Democratic or Republican candidates, I include in table 1.5 only the top quartile (for Democrat) and bottom quartile (for Republican) by number of links of the friendship counties. I call these variables *friend_dd_exp_{ct}* and *friend_rr_exp_{ct}*. They are defined as follows:

$$\begin{aligned}
 \text{friend_dd_exp}_{ct} = \sum_{j=1}^{N-1} & \left(\frac{SCI_{c,j}}{\sum_{k=1}^{N-1} SCI_{c,k}} \times (\text{Vote margin}_{jt} - \text{Vote margin}_{j,t-4}) \right. \\
 & \left. \times \mathbb{I}\{\text{Vote margin}_{jt} - \text{Vote margin}_{j,t-4} > q_{0.75}\} \right)
 \end{aligned} \tag{1.10}$$

¹⁴For the first stage results of the 2SLS analysis please consult table 1.8.

$$\begin{aligned}
friend_rr_exp_{ct} = \sum_{j=1}^{N-1} & \left(\frac{SCI_{c,j}}{\sum_{k=1}^{N-1} SCI_{c,k}} \times (Vote\ margin_{jt} - Vote\ margin_{j,t-4}) \right. \\
& \left. \times \mathbb{I}\{Vote\ margin_{jt} - Vote\ margin_{j,t-4} < q_{0.25}\} \right)
\end{aligned} \tag{1.11}$$

The OLS results look very similar to the model where Democrat and Republican exposure variables are split along the median. In the 2SLS specification similarly there is no substantive change from the median split model.¹⁵ It seems that the above results are not driven by counties that swung extremely into either political direction.

Table 1.7 depicts the results for influence of friendship county exposure on Republican-leaning and Democrat-leaning counties separately. Similar to Dahl et al. (2014)'s strategy, I analyze the effect of friendship influence on predominantly Democrat counties — only those counties are examined where Democrat-Republican vote margin in each year is above the median value of -33%. In half of the counties Democrat margin over Republicans was -33% or higher.¹⁶ The first column of this table denotes the greater than -.33 group - mostly Democratic counties. As is apparent in the table for 2012 in the mostly Republican counties the coefficient of influence of other Democrat counties (weighed by the SCI) is 23.33. This is quite a bit lower than the influence coefficient in any other 2SLS specification above.

In the above median specification the coefficient of influence of Democratic counties is statistically insignificant for 2012. It appears that there was coalition building among already-Republican counties but in predominantly Democratic counties the effect of ties with Democrat counties was not perceptible, at least in this model. This specification may shine light on the positive coefficient on *friend_rr_exp* (the effect of Republican county connections). In the above-median (>-33%) equation the coefficient is 11.52 and significant at the 1% level. However in the below-median equation the magnitude falls to 3.03. Statistical significance remains at the 1% level. The positive coefficient of the *friend_rr_exp*, which

¹⁵For the first stage results of the 2SLS analysis please consult table 1.8.

¹⁶Counties where the (*Democrat%* – *Republican%*) margin is below the median, i.e. Republican favored counties, outweigh Democrat favored counties by two to one These are largely rural counties with low population densities.

is the effect of Republican connections on vote margin, was driven by mostly Democratic counties in 2012. This coefficient should be negative if the effect of Republican ties reduces the (*Democrat – Republican*) margin. In 2020 and 2016 the coefficient on *friend_r_exp* and *friend_rr_exp* is negative where statistical significance is attained. In the above median (>33) counties the effect of *friend_rr_exp* is statistically insignificant in 2016. It is significant at the 5% level and negative for the below median (Republican) counties for both 2016 and 2020.

As for *friend_dd_exp* in year 2016, the coefficient of influence on below median vote margin counties (ie strongly Republican counties) is higher at 14.56% than the influence on above median vote margin counties (Democrat), 13.35%. It is possible that the formation of coalitions on the Democratic side is concluding. This is consistent with the decrease of the *friend_exp* coefficient from election to election in the level regressions discussed in table 1.2.

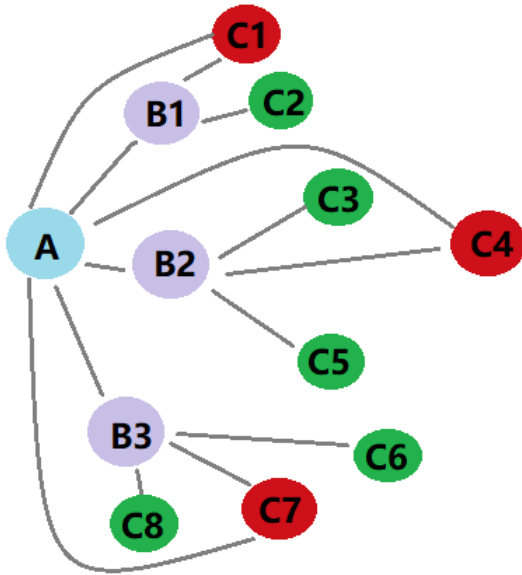
The stronger effect (coefficient 18.2% in 2016 and 14.56% in 2020) of *friend_dd_exp* on Republican counties than Democrat counties (9.7% in 2016, 13.35% in 2020) indicates greater mobilization in these counties by Democrats. When considered in conjunction with the effect of *friend_rr_exp* effect on below median (Republican) counties this also indicates strong political polarization in below-median vote margin (Republican) counties. The effect of *friend_rr_exp* on above median (Democrat) is positive in all years (insignificant coefficient in 2016).

1.6 Conclusion

This paper looks at the effect of social networks on decision making. More specifically it looks at the effect of social networks on presidential vote decisions. By utilizing county-to-county social connectedness data from Facebook corporation I discern an approximately 1 percentage point increase in vote margin for a 1 percentage point increase in the social network vote margin for all years examined, controlling for the usual vote determinants in

addition to time and state fixed effects. To alleviate concerns about endogeneity stemming from simultaneous action I implement a 2SLS approach. This approach confirms the existence of social network effects on county vote decisions. I also look at treatment intensity by breaking up the social network into the top and bottom quartiles of vote margin. While more extreme Republican or Democrat counties are not driving the social network effects I do find heterogeneity of outcomes when the dependent variable is split up into Democrat- and Republican-favoring counties. The effect of the Democrat part of the social network is positive in both Republican and Democrat focal counties; however the influence of the Republican portion of a county's social network is positive for above median counties (Democrat) but negative for Republican ones. Furthermore this wedge of the effect of the Republican portion of the social network on above (Democrat) and below (Republican) median focal counties grows overtime. This points to increasing political polarization.

Figure 1.1: Intransitive Triads



In the figure focal county *A* has friendship links to the lavender and red colored counties. *A* has no direct links to the green counties, only second order links through the lavender ones. Therefore the characteristics of the green counties are optimal to construct an instrumental variable that meets the exogeneity restriction. The red counties are directly linked to the focal county, in addition to second order connections through the lavender ones. The red counties can not be used to construct the instrumental variable.

1.7 Appendix

Table 1.1: Summary statistics

	All years			2008			2012			2016			2020		
	Mean	Std.Dev		Mean	Std.Dev		Mean	Std.Dev		Mean	Std.Dev		Mean	Std.Dev	
vote margin	-0.25062	0.308093		-0.15338	0.276186		-0.21232	0.295061		-0.31837	0.30681		-0.31843	0.31967	
friend_exp	-0.07874	0.152542		0.008408	0.086383		-0.08635	0.149583		-0.13035	0.159694		-0.10668	0.163276	
friend_exp_outstate	-0.02715	0.101978		0.008408	0.086383		-0.0308	0.094245		-0.05643	0.106733		-0.0298	0.108304	
friend_d_exp	0.020489	0.018081					0.004223	0.005213		0.023355	0.019862		0.033888	0.010334	
friend_r_exp	-0.04249	0.036114					-0.04991	0.027489		-0.06736	0.036963		-0.01022	0.008729	
percent of population over 65	0.169889	0.045738		0.15357	0.042026		0.161014	0.042093		0.176383	0.044408		0.188588	0.046072	
percent of population Hispanic	0.085813	0.134521		0.075526	0.128149		0.083342	0.133151		0.089848	0.136586		0.094537	0.139247	
percent of population Black	0.090247	0.145193		0.08903	0.144486		0.090178	0.145828		0.090608	0.145283		0.09117	0.145234	
percent of population with a bachelor degree who is over 25	0.221525	0.096717		0.186691	0.084928		0.27223	0.093813		0.207674	0.091131		0.219505	0.095527	
unemployment rate	0.069492	0.03381		0.069045	0.030858		0.086157	0.037142		0.070261	0.032034		0.052505	0.025465	
population of county	100695.5	321697.6		96840.61	311027		99303.71	315447.2		102330.3	326970.3		104307.5	332972.5	
per capita income	24607.66	6269.502		22123.06	5250.619		23316.67	5448.029		24965.54	5924.709		28024.57	6727.634	
N		12420			3105			3105			3105			3105	

¹ Observations at the county level. Vote data was obtained from MIT Election Lab. All other variables come from the American Community Survey.

² The logarithm of population and income are used in subsequent analysis.

Table 1.2: Level Regressions

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
	All	2008	2012	2016	2020	All	2008	2012	2016	2020
friend_exp	1.044*** (0.0169)	1.508*** (0.0584)	1.317*** (0.0376)	0.982*** (0.0326)	0.954*** (0.0313)	1.394*** (0.0250)	1.508*** (0.0584)	1.482*** (0.0547)	1.090*** (0.0453)	1.079*** (0.0430)
friend_exp_outstate						0.361*** (0.0374)	0.403*** (0.0786)	0.293*** (0.0817)	0.281*** (0.0679)	0.330*** (0.0639)
percent population over 65	0.343*** (0.0366)	0.403*** (0.0786)	0.349*** (0.0769)	0.310*** (0.0650)	0.370*** (0.0614)	0.566*** (0.0374)	0.486*** (0.0786)	0.597*** (0.0817)	0.696*** (0.0679)	0.523*** (0.0639)
percent population Hispanic	0.528*** (0.0148)	0.486*** (0.0313)	0.458*** (0.0301)	0.593*** (0.0269)	0.466*** (0.0256)	0.842*** (0.0151)	0.486*** (0.0313)	0.597*** (0.0312)	0.696*** (0.0273)	0.523*** (0.0264)
percent population black	0.877*** (0.0151)	0.717*** (0.0322)	0.732*** (0.0307)	0.924*** (0.0286)	0.945*** (0.0279)	0.842*** (0.0160)	0.717*** (0.0322)	0.811*** (0.0330)	0.985*** (0.0303)	1.010*** (0.0291)
percent population over 25 with bachelor degree	1.039*** (0.0249)	0.641*** (0.0557)	0.741*** (0.0514)	1.288*** (0.0444)	1.429*** (0.0423)	0.938*** (0.0262)	0.641*** (0.0557)	0.641*** (0.0560)	1.206*** (0.0482)	1.324*** (0.0459)
unemployment rate	1.233*** (0.0552)	1.475*** (0.119)	1.056*** (0.104)	1.057*** (0.103)	1.269*** (0.114)	1.398*** (0.0562)	1.475*** (0.119)	1.421*** (0.110)	1.429*** (0.106)	1.507*** (0.118)
log of population	0.00669*** (0.00137)	0.00267 (0.00279)	-0.00283 (0.00274)	0.00512** (0.00247)	0.0139*** (0.00246)	0.00666*** (0.00140)	0.00267 (0.00279)	-0.00182 (0.00293)	0.00993*** (0.00257)	0.0204*** (0.00252)
log of per capita income	-0.260*** (0.0108)	-0.219*** (0.0226)	-0.316*** (0.0224)	-0.294*** (0.0198)	-0.291*** (0.0197)	-0.233*** (0.0110)	-0.219*** (0.0226)	-0.237*** (0.0236)	-0.214*** (0.0204)	-0.214*** (0.0201)
N	12419	3104	3105	3105	3105	12419	3104	3105	3105	3105

¹ Column one and 6 report results from panel regression of years 2008, 2012, 2016, 2020. Time and State fixed effects are used. Columns 2 through 5 annual regressions of the comprehensive network are shown; state fixed effects are used. Columns 7 through 10 only the effects of the out-of-state network are looked at; state fixed affects are used. Observations are at the county level. Vote data was obtained from MIT Election Lab. County-to-county social connections come from Facebook. All other variables come from the American Community Survey.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 1.3: Vote margin shift regressions. All connected counties.

	(1)	(2)	(3)	(4)
	All	2012	2016	2020
friend_d_exp	1.988*** (0.150)	7.412*** (1.165)	0.530* (0.296)	4.058*** (0.358)
friend_r_exp	-0.408*** (0.0805)	0.840*** (0.211)	-1.149*** (0.157)	-3.406*** (0.411)
percent population over 65	0.429*** (0.0463)	0.366*** (0.0902)	0.280*** (0.0741)	0.433*** (0.0689)
percent population Hispanic	0.816*** (0.0179)	0.737*** (0.0347)	0.926*** (0.0287)	0.610*** (0.0316)
percent population black	1.382*** (0.0170)	1.142*** (0.0355)	1.459*** (0.0283)	1.427*** (0.0257)
percent population with a bachelor degree over 25	1.496*** (0.0303)	1.053*** (0.0591)	1.700*** (0.0498)	1.685*** (0.0461)
unemployment rate	1.564*** (0.0693)	1.689*** (0.121)	1.582*** (0.114)	1.541*** (0.127)
log of population	0.0273*** (0.00170)	0.0186*** (0.00315)	0.0303*** (0.00274)	0.0355*** (0.00263)
log of per capita income	-0.193*** (0.0138)	-0.198*** (0.0261)	-0.162*** (0.0222)	-0.197*** (0.0217)
N	9315	3105	3105	3105

¹ Friend.d_exp is the county level summation of 4 year (Democrat - Republican) vote margin changes where the change over 4 years was above the median, weighted by the SCI social connection data. Friend.r_exp is the county level summation of the changes in the 4 year vote margin where these changes were below the median. Also weighted by the SCI at the county level. Observations are at the county level. Vote data was obtained from MIT Election Lab. County-to-county social connections come from Facebook. All other variables come from the American Community Survey.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 1.4: Restricted network OLS and 2SLS shift regressions

	OLS				2SLS			
	(1) All	(2) 2012	(3) 2016	(4) 2020	(5) All	(6) 2012	(7) 2016	(8) 2020
friend_d_exp	1.960*** (0.148)	7.351*** (1.156)	0.523* (0.292)	3.959*** (0.352)	32.01*** (3.936)	219.1*** (50.72)	37.14*** (6.012)	22.76*** (1.925)
friend_r_exp	-0.418*** (0.0794)	0.805*** (0.209)	-1.172*** (0.154)	-3.359*** (0.407)	-9.256*** (1.896)	3.633 (3.487)	-6.289*** (1.954)	-18.35*** (2.920)
percent population over 65	0.421*** (0.0464)	0.359*** (0.0903)	0.274*** (0.0741)	0.417*** (0.0691)	0.623*** (0.113)	0.739** (0.335)	1.078*** (0.216)	0.645*** (0.102)
percent of population Hispanic	0.809*** (0.0179)	0.732*** (0.0349)	0.918*** (0.0287)	0.604*** (0.0317)	0.797*** (0.0428)	-1.013** (0.394)	0.707*** (0.0820)	0.179 (0.116)
percent of population black	1.382*** (0.0170)	1.144*** (0.0354)	1.459*** (0.0281)	1.426*** (0.0257)	1.489*** (0.0577)	-2.015*** (0.703)	1.517*** (0.106)	1.441*** (0.0364)
percent of population over 25 with a bachelor degree	1.489*** (0.0305)	1.045*** (0.0595)	1.687*** (0.0500)	1.683*** (0.0463)	1.219*** (0.0788)	0.469* (0.246)	0.827*** (0.175)	1.413*** (0.0740)
unemployment rate	1.587*** (0.0694)	1.702*** (0.121)	1.615*** (0.114)	1.569*** (0.128)	1.684*** (0.175)	-0.535 (0.639)	1.815*** (0.284)	1.141*** (0.198)
log of population	0.0274*** (0.00171)	0.0188*** (0.00317)	0.0309*** (0.00275)	0.0352*** (0.00266)	0.00660 (0.00469)	0.0720*** (0.0173)	-0.0269** (0.0107)	0.0223*** (0.00419)
log of per capita income	-0.188*** (0.0138)	-0.194*** (0.0262)	-0.153*** (0.0222)	-0.193*** (0.0218)	-0.182*** (0.0368)	-0.584*** (0.128)	-0.182*** (0.0604)	-0.293*** (0.0328)
N	9306	3102	3102	3102	9306	3102	3102	3102

¹ friend_d_exp is the county level summation of 4 year (Democrat - Republican) vote margin changes where the change over 4 years was above the median, weighted by the SCI social connection data. friend_r_exp is the county level summation of the changes in the 4 year vote margin where these changes were below the median. Also weighted by the SCI at the county level. The restricted network data set is used: only above median SCI connections are counted to better discern Intransitive Triads. In columns 1 through 4 OLS results are reported for the restricted network. In columns 5 through 8 the 2SLS results are reported. For the 2SLS results the coefficients represent the Local Average Treatment Effects. The effect of the network for compliers only. Observations are at the county level. Vote data was obtained from MIT Election Lab. County-to-county social connections come from Facebook. All other variables come from the American Community Survey.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 1.5: Restricted network top and bottom quartiles shift OLS and 2SLS regressions

	OLS				2SLS			
	(1) All	(2) 2012	(3) 2016	(4) 2020	(5) All	(6) 2012	(7) 2016	(8) 2020
friend_d_exp top quartile	1.587*** (0.136)	7.755*** (1.108)	0.0817 (0.280)	3.094*** (0.326)	9.919*** (1.607)	151.8*** (30.200)	33.24*** (5.300)	22.65*** (2.186)
friend_r_exp bottom quartile	-0.350*** (0.078)	0.767*** (0.181)	-1.165*** (0.134)	-2.716*** (0.393)	2.654** (1.040)	3.345 (2.230)	-6.748*** (2.234)	-15.34*** (3.039)
percent population over 65yo	0.433*** (0.046)	0.360*** (0.090)	0.286*** (0.074)	0.439*** (0.069)	0.659*** (0.067)	0.592** (0.236)	1.044*** (0.210)	0.811*** (0.115)
percent population Hispanic	0.813*** (0.018)	0.749*** (0.034)	0.911*** (0.028)	0.619*** (0.032)	0.782*** (0.025)	-0.333 (0.219)	0.638*** (0.079)	0.244* (0.126)
percent population black	1.379*** (0.017)	1.151*** (0.035)	1.449*** (0.028)	1.427*** (0.026)	1.268*** (0.029)	-1.084** (0.448)	1.505*** (0.103)	1.443*** (0.039)
percent population over 25 yo with a bachelor degree	1.503*** (0.030)	1.058*** (0.059)	1.687*** (0.049)	1.717*** (0.046)	1.335*** (0.044)	0.741*** (0.162)	0.790*** (0.176)	1.582*** (0.075)
unemployment rate	1.563*** (0.069)	1.692*** (0.121)	1.552*** (0.114)	1.593*** (0.128)	1.378*** (0.098)	0.237 (0.408)	1.650*** (0.278)	1.474*** (0.206)
log of population	0.0276*** (0.002)	0.0180*** (0.003)	0.0305*** (0.003)	0.0360*** (0.003)	0.0160*** (0.003)	0.0453*** (0.010)	-0.0222** (0.010)	0.0194*** (0.005)
log of per capita income	-0.193*** (0.014)	-0.197*** (0.026)	-0.170*** (0.022)	-0.192*** (0.022)	-0.261*** (0.021)	-0.446*** (0.084)	-0.208*** (0.056)	-0.302*** (0.036)
N	9315	3105	3105	3105	9306	3102	3102	3102

¹ Friend_d_exp is the county level summation of 4 year (Democrat - Republican) vote margin changes where the change over 4 years was above the top quartile, weighted by the SCI social connection data. Friend_r_exp is the county level summation of the changes in the 4 year vote margin where these changes were below the bottom quartile. Also weighted by the SCI at the county level. The restricted network data set is used: only above median SCI connections are counted to better discern Intransitive Triads. In columns 1 through 4 OLS results are reported for the restricted network. In columns 5 through 8 the 2SLS results are reported. For the 2SLS results the coefficients represent the Local Average Treatment Effects. The effect of the network for compliers only. Observations are at the county level. Vote data was obtained from MIT Election Lab. County-to-county social connections come from Facebook. All other variables come from the American Community Survey.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 1.6: Split dependent variable along median 2SLS analysis

2SLS	All		2012		2016		2020	
	(1) >= Median	(2) <Median	(3) >= Median	(4) <Median	(5) >= Median	(6) <Median	(7) >= Median	(8) <Median
friend_d_exp	3.111*** (0.239)	-1.750*** (0.191)	13.59 (34.92)	25.16** (11.54)	14.37*** (4.513)	20.07*** (4.136)	14.77*** (3.187)	14.89*** (1.966)
friend_r_exp	0.0163 (0.134)	-0.0170 (0.0650)	19.78*** (7.130)	3.011*** (0.860)	-0.178 (1.836)	-2.868** (1.213)	3.719 (3.212)	-21.73*** (3.781)
percent of population over 65	0.0730 (0.0653)	0.523*** (0.0419)	0.492 (0.345)	0.576*** (0.102)	0.442*** (0.159)	1.028*** (0.199)	0.527*** (0.152)	0.793*** (0.104)
percent of population Hispanic	0.702*** (0.0248)	0.381*** (0.0192)	0.0790 (0.303)	0.0940 (0.0663)	0.736*** (0.0651)	0.387*** (0.0865)	1.026*** (0.139)	-0.0321 (0.0901)
percent of population black	1.214*** (0.0233)	1.004*** (0.0294)	0.234 (0.490)	0.362** (0.175)	1.258*** (0.0906)	1.220*** (0.140)	1.517*** (0.0554)	1.224*** (0.0745)
percent of population over 25 with bachelor degree	1.502*** (0.0381)	0.326*** (0.0320)	1.038*** (0.137)	-0.0824 (0.0772)	1.274*** (0.108)	-0.420* (0.218)	1.819*** (0.0835)	0.170* (0.0890)
unemployment rate	1.136*** (0.0893)	0.606*** (0.0667)	0.0736 (0.482)	0.625*** (0.158)	1.090*** (0.191)	0.830*** (0.274)	1.484*** (0.273)	0.549*** (0.172)
log of population	-0.0110*** (0.00231)	0.0470*** (0.00175)	-0.0171 (0.0167)	0.0391*** (0.00431)	-0.0299*** (0.00661)	0.000128 (0.0112)	-0.00796 (0.00540)	0.0417*** (0.00552)
log of per capita income	-0.238*** (0.0192)	0.0414*** (0.0130)	-0.429*** (0.116)	-0.0978** (0.0407)	-0.281*** (0.0416)	0.0978* (0.0542)	-0.306*** (0.0439)	-0.0481 (0.0342)
N	4653	4653	1551	1551	1551	1551	1551	1551

¹ friend_d_exp is the county level summation of 4 year (Democrat - Republican) vote margin changes where the change over 4 years was above the median, weighted by the SCI social connection data. friend_r_exp is the county level summation of the changes in the 4 year vote margin where these changes were below the median, also weighted by the SCI at the county level. The restricted network data set is used: only above median SCI connections are counted to better discern Intransitive Triads. This table only reports 2SLS derived results. In odd numbered columns the dependent variable is greater than or equal to the median vote margin. In even numbered columns the dependent variable is below median vote margins. In the 2SLS results the coefficients represent the Local Average Treatment Effects. The effect of the network is for compliers only. Observations are at the county level. Vote data was obtained from MIT Election Lab. County-to-county social connections come from Facebook. All other variables come from the American Community Survey.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 1.7: Split dependent variable 2SLS results for top and bottom quartile independent variables of interest

2SLS	2012				2016				2020			
	(1) >= Median	(2) <Median	(3) >= Median	(4) <Median	(5) >= Median	(6) <Median	(7) >= Median	(8) <Median				
friend_dd_exp	3.754*** (0.288)	-1.610*** (0.201)	28.54 (19.22)	23.33** (9.968)	9.702*** (3.634)	18.20*** (3.600)	13.35*** (3.971)	14.56*** (1.985)				
friend_rr_exp	1.843*** (0.350)	-0.549*** (0.125)	11.52*** (3.306)	3.038*** (0.761)	3.052 (1.959)	-2.972*** (1.124)	7.642** (3.227)	-20.04*** (3.814)				
percent of population over 65	0.119* (0.0694)	0.503*** (0.0420)	0.335 (0.232)	0.565*** (0.101)	0.414*** (0.156)	0.987*** (0.189)	0.599*** (0.188)	0.903*** (0.114)				
percent of population Hispanic	0.699*** (0.0262)	0.385*** (0.0193)	0.270 (0.171)	0.134** (0.0576)	0.700*** (0.0553)	0.334*** (0.0781)	1.154*** (0.151)	-0.0290 (0.0967)				
percent of population black	1.172*** (0.0261)	1.015*** (0.0300)	0.496* (0.296)	0.484*** (0.165)	1.147*** (0.0862)	1.234*** (0.136)	1.500*** (0.0612)	1.246*** (0.0798)				
percent of population over 25 with bachelor degree	1.460*** (0.0401)	0.349*** (0.0319)	1.064*** (0.106)	-0.0748 (0.0752)	1.279*** (0.109)	-0.441** (0.214)	1.962*** (0.0879)	0.225** (0.0919)				
unemployment rate	1.084*** (0.0936)	0.624*** (0.0671)	0.273 (0.316)	0.631*** (0.154)	1.152*** (0.194)	0.803*** (0.264)	1.760*** (0.277)	0.739*** (0.192)				
log of population	-0.0140*** (0.00249)	0.0472*** (0.00178)	-0.0115 (0.00867)	0.0366*** (0.00401)	-0.0287*** (0.00643)	0.00301 (0.0107)	-0.0127** (0.00583)	0.0425*** (0.00616)				
log of per capita income	-0.250*** (0.0207)	0.0435*** (0.0135)	-0.390*** (0.0825)	-0.0724** (0.0361)	-0.282*** (0.0402)	0.0723 (0.0488)	-0.302*** (0.0486)	-0.0407 (0.0363)				
N	4652	4654	1550	1552	1551	1551	1551	1551				

¹ Friend_d_exp is the county level summation of 4 year (Democrat - Republican) vote margin changes where the change over 4 years was above the top quartile, weighted by the SCI social connection data. Friend_r_exp is the county level summation of the changes in the 4 year vote margin where these changes were below the bottom quartile, also weighted by the SCI at the county level. The restricted network data set is used: only above median SCI connections are counted to better discern Intransitive Triads. This table only reports 2SLS derived results. In odd numbered columns the dependent variable is greater than or equal to the median vote margin. In even numbered columns the dependent variable is below median vote margins. The 2SLS results represent the Local Average Treatment Effects. The effect of the network for compliers only. Observations are at the county level. Vote data was obtained from MIT Election Lab. County-to-county social connections come from Facebook. All other variables come from the American Community Survey.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 1.8: First stage of the 2SLS restricted network shift regressions

	all			2012			2016			2020		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	
	friend_d.exp	friend_r.exp	friend_d.exp	friend_r.exp	friend_d.exp	friend_r.exp	friend_d.exp	friend_r.exp	friend_d.exp	friend_r.exp	friend_d.exp	friend_r.exp
IT average of percent population over 25 with bachelor degree	0.0734 (0.0541)	-0.233** (0.100)	0.00628 (0.0240)	-0.812*** (0.130)	-0.268*** (0.0821)	-1.470*** (0.151)	-0.578*** (0.0655)	-0.307*** (0.0581)				
IT average of percent of population over 65	-0.144 (0.114)	0.247 (0.211)	-0.0158 (0.0555)	-0.731** (0.301)	-0.250 (0.176)	-1.878*** (0.324)	0.162 (0.128)	-0.0973 (0.113)				
IT average of percentage population black	0.0156 (0.0227)	-0.0231 (0.0421)	-0.00712 (0.00917)	-0.239*** (0.0496)	-0.0457 (0.0355)	-0.0772 (0.0653)	0.184*** (0.0275)	0.0828*** (0.0244)				
IT average of percent population Hispanic	-0.109*** (0.0295)	-0.341*** (0.0548)	0.0153 (0.0114)	-0.204*** (0.0620)	-0.0967** (0.0384)	-0.646*** (0.0707)	-0.0137 (0.0309)	-0.253*** (0.0274)				
IT average of unemployment rate	-1.202*** (0.119)	-2.665*** (0.221)	-0.0854 (0.0606)	0.702** (0.328)	-0.307 (0.295)	-1.551*** (0.543)	-1.569*** (0.288)	-0.671*** (0.255)				
IT average of log of population	0.00436 (0.00305)	0.0250*** (0.00565)	0.0000363 (0.00142)	-0.0375*** (0.00770)	-0.00625 (0.00466)	-0.00107 (0.00858)	0.00776** (0.00357)	0.00476 (0.00317)				
IT average of log of per_capita_income	-0.227*** (0.0287)	-0.320*** (0.0533)	-0.0264* (0.0137)	0.247*** (0.0740)	-0.0530 (0.0489)	0.199** (0.0900)	0.00580 (0.0395)	0.0401 (0.0351)				
percent of population over 65	-0.0218*** (0.00383)	-0.0304*** (0.00710)	-0.00257* (0.00148)	-0.0259*** (0.00802)	-0.0354*** (0.00490)	-0.0590*** (0.00901)	-0.0284*** (0.00384)	-0.0102*** (0.00340)				
percent of population Hispanic	-0.000453 (0.00172)	-0.00682** (0.00320)	0.00825*** (0.000647)	0.0251*** (0.00350)	0.00625*** (0.00220)	0.0110*** (0.00404)	-0.0136*** (0.00178)	-0.0514*** (0.00158)				
percent of population black	0.00867*** (0.00221)	0.0220*** (0.00411)	0.0142*** (0.000834)	0.0272*** (0.00451)	0.00354 (0.00285)	0.0381*** (0.00524)	0.00835*** (0.00226)	0.00178 (0.00201)				
percent of population over 25 with a bachelor degree	0.00902*** (0.00305)	-0.000984 (0.00567)	0.00150 (0.00118)	-0.0215*** (0.00636)	0.0137*** (0.00397)	0.00780 (0.00731)	-0.0189*** (0.00312)	-0.0200*** (0.00277)				
unemployment rate	-0.00421 (0.00569)	0.0237** (0.0106)	0.00850*** (0.00196)	0.0559*** (0.0106)	-0.00753 (0.00741)	-0.00163 (0.0136)	0.00317 (0.00681)	-0.0274*** (0.00605)				
log of population	0.000884*** (0.000165)	0.000869*** (0.000305)	-0.000325*** (0.0000615)	-0.00138*** (0.000333)	0.00144*** (0.000211)	0.00290*** (0.000388)	0.000895*** (0.000169)	0.00103*** (0.000150)				
log of per capita income	0.00139 (0.00114)	0.0115*** (0.00211)	0.00124*** (0.000432)	0.0186*** (0.00234)	0.00274* (0.00145)	0.0214*** (0.00267)	0.00878*** (0.00118)	0.00432*** (0.00105)				
CLR statistic	1129.450		381.644			374.459		435.309				
CLR p-value	0.000		0.000			0.000		0.000				
Wald statistic	82.070		24.525			43.566		146.479				
Wald p-value	0.000		0.000			0.000		0.000				
N	9306	9306	3102	3102	3102	3102	3102	3102				3102

¹ The instrument is obtained from the presence of Intransitive Triads (IT). The exogenous variables of county friends of friends who are not friends of the focal county. Weak instrument test statistics are reported.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 1.9: First stage results for 2SLS restricted network regressions with top and bottom quartiles

	all													
	2012				2016				2020					
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)
	friend_dd_exp	friend_rr_exp	friend_dd_exp	friend_rr_exp	friend_dd_exp	friend_rr_exp	friend_dd_exp	friend_rr_exp	friend_dd_exp	friend_rr_exp	friend_dd_exp	friend_rr_exp	friend_dd_exp	friend_rr_exp
I T average of percentage over	0.146***	-0.535***	-0.002	-1.145***	-0.276***	-1.164***	-0.476***	-0.281***						
25 with bachelor degree	(0.055)	(0.096)	(0.025)	(0.147)	(0.083)	(0.172)	(0.069)	(0.058)						
I T average percent of population over 65	0.210*	-0.203	-0.048	-0.790**	-0.244	-2.030***	0.357***	-0.112						
	(0.116)	(0.203)	(0.057)	(0.340)	(0.177)	(0.368)	(0.135)	(0.113)						
I T average population percent black	0.0328	-0.022	-0.0211**	-0.160***	-0.044	-0.072	0.171***	0.0729***						
	(0.023)	(0.040)	(0.009)	(0.056)	(0.036)	(0.074)	(0.029)	(0.024)						
I T average population percent Hispanic	-0.0838***	-0.289***	0.009	-0.230***	-0.0994**	-0.535***	0.0037	-0.247***						
	(0.030)	(0.053)	(0.012)	(0.070)	(0.039)	(0.080)	(0.033)	(0.027)						
I T average unemployment rate	-1.211***	-1.388***	-0.104*	0.706*	-0.36	-0.819	-1.044***	-0.556**						
	(0.121)	(0.212)	(0.062)	(0.372)	(0.297)	(0.616)	(0.304)	(0.255)						
IT average log of population	0.00674**	0.0141***	-0.001	-0.0269***	-0.005	-0.008	0.004	0.0043						
	(0.003)	(0.005)	(0.001)	(0.009)	(0.005)	(0.010)	(0.004)	(0.003)						
I T average log per_capita_income	-0.238***	-0.066	-0.022	0.355***	-0.054	0.229**	0.033	0.0477						
	(0.029)	(0.051)	(0.014)	(0.084)	(0.049)	(0.102)	(0.042)	(0.035)						
percentage population over 65yo	-0.0222***	-0.0246***	-0.002	-0.0215**	-0.0349***	-0.0421***	-0.0313***	-0.0103***						
	(0.004)	(0.007)	(0.002)	(0.009)	(0.005)	(0.010)	(0.004)	(0.003)						
percent population Hispanic	0.000201	-0.0134***	0.00722***	0.0149***	0.00640***	-2E-04	-0.0107***	-0.0514***						
	(0.002)	(0.003)	(0.001)	(0.004)	(0.002)	(0.005)	(0.002)	(0.002)						
percent population black	0.00881***	0.0134***	0.0136***	0.0103**	0.0029	0.0325***	0.00808***	0.0011						
	(0.002)	(0.004)	(0.001)	(0.005)	(0.003)	(0.006)	(0.002)	(0.002)						
percent population over	0.00889***	-0.0128**	0.0011	-0.0327***	0.0144***	0.0008	-0.0207***	-0.0194***						
25 with bachelor degree	(0.003)	(0.005)	(0.001)	(0.007)	(0.004)	(0.008)	(0.003)	(0.003)						
unemployment rate	-0.00582	0.0248**	0.00729***	0.0633***	-0.008	-0.025	-0.008	-0.0284***						
	(0.006)	(0.010)	(0.002)	(0.012)	(0.007)	(0.016)	(0.007)	(0.006)						
log of population	0.000816***	0.000883***	-0.000335***	-0.000731*	0.00143***	0.00249***	0.000835***	0.00101***						
	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)						
log of per capita income	0.00124	0.0107***	0.00106**	0.0198***	0.00243*	0.0126***	0.00857***	0.00404***						
	(0.001)	(0.002)	(0.000)	(0.003)	(0.001)	(0.003)	(0.001)	(0.001)						
CLR statistic	1092.625		373.584		375.270		425.215							
CLR p-value	0.000		0.000		0.000		0.000							
Wald statistic	115.215		34.502		42.924		109.565							
Wald p-value	0.000		0.000		0.000		0.000							
N	9306	9306	3102	3102	3102	3102	3102	3102						

¹ The instrument is obtained from the presence of Intransitive Triads (IT). The exogenous variables of county friends of friends who are not friends of the focal county. Weak instrument test statistics are reported.

² Standard errors in parentheses, * p < 0.1, ** p < 0.05, *** p < 0.01.

Chapter 2

Incorporating Bridging Social Capital into Social Capital Measures

2.1 Introduction

For geographically comprehensive studies researchers measure social capital using rates of membership in social organizations. However, membership numbers may be picking up mostly one type of social capital — bonding. As such, these measures suffer from missing an important dimension of social capital — bridging social capital. Why is this a problem? Can it be ameliorated?

While social capital is typically thought to be valuable for a society, some recent studies point out the potentially harmful effects of social capital. Satyanath et al. (2017) note that those municipalities in pre-World War II Germany with the highest level of associational membership rates were quicker to be taken over by the Nazi party. Social organizations helped spread Nazi ideology. The measure Satyanath et al. (2017) use is very similar to the measures being used for current social capital studies. The reason these measures seem to find harmful effects of social capital is because they miss the bridging element of it. By its very nature, bridging social capital is more difficult to identify. The building blocks of

bridging social capital are weak links between people. These are links between acquaintances — people who barely know each other. Communities with a high level of bonding and a low level of bridging social capital are insular, more mistrustful of others, and have difficulties in achieving a society-wide consensus (Bobo (1988)). Woolcock (1998) notes communities with high levels of bonding social capital get by and those with bridging social capital get ahead. Therefore, measures that pick up bonding social capital will suffer from mis-measurement problems.

This study uses recently available county-to-county social connectedness data from Facebook to remedy the social capital mis-measurement problem. The research design rests on the idea that geographically distant ties contain more weak ties, and therefore higher levels of bridging social capital. These distant ties, combined with the current, within-county measures of social capital constitute a more comprehensive measure of social capital. The new measure picks up both bonding and bridging social capital. As such, the measure safeguards against false labeling of insular communities as high-social-capital when in fact they have high level of one variety of social capital. In this manner the new measure confirms the beneficial qualities of social capital.

The manner in which bridging and bonding social capital are combined into a comprehensive social capital variable uses social network topology. The Social Connectedness Index (SCI) dataset provides a normalized county-to-county measure of social connections as measured by Facebook connections between the population of U.S. county pairs. Since Facebook does not provide individual-level connectivity data, within-county connections provide very little information. However, under the assumption that social ties between counties contain higher levels of bridging social capital, I can use the SCI to measure the cross-county variation in bridging social capital by using eigenvector centrality. To combine this measure of bridging social capital with bonding social capital, I use the closely-related alpha centrality method. Alpha centrality allows me to combine a network-derived centrality measure with an exogenous (to the network) vector of county-level bonding social capital. Furthermore,

the new variable has a tunable parameter, alpha. I can vary alpha to look for social links at the local or global network level. This allows me to use grid search, a machine learning tool, to tune the value of alpha. By varying the level of alpha, grid search picks out the one that minimizes mean squared errors in a variety of regressions I run. Alpha centrality also allows me to search for counties that have the largest level of bonding social capital. For these reasons, the new variable is an improvement over the old ones that measure local-level associational membership rates, missing bridging social capital.

In the empirical section of this work I compare the new, comprehensive measure of social capital to some of the previously used ones. In one instance, I replicate the Goetz et al. (2012) study of hate-group formation for a different year using their social capital variable and the new, comprehensive one. The social capital variable Goetz et al. (2012) provide has the expected negative sign and is significant in their study that uses data from 2008, confirming that higher levels of social capital inhibit hate-group formation. However, when their empirical study is replicated using variables from 2005, the coefficient on their social capital variable becomes insignificant. When the new, comprehensive social capital measure is substituted for their measure, the coefficient on this variable is negative and significant at the 1% level. Including bridging social capital reduces mis-measurement problems in the social capital variable.

I also deploy the new social capital variable in estimating the impact of social capital on crime levels, unemployment and obesity rates. Combined social capital is positive and significant at the 1% level for total crime and violent crime, significant at the 5% level for murder, but not significant in the property crime model. Bonding social capital is statistically insignificant in all regressions.

For the county-level obesity and unemployment models the level of the alpha parameter that minimizes mean square errors is zero. At alpha equaling zero, the combined social capital measure replicates bonding social capital; therefore, the equations are identical for bonding and combined social capital. The coefficient on social capital is the expected negative sign

and significant at the 1% level in the unemployment and at the 5% level in the obesity equations.

To delve deeper into the puzzling estimated positive effect of social capital on crime, I examine the effect of increasing the radius of assessed network centrality. Alpha-centrality-derived measures have the inherent benefit of allowing one to conduct local vs. global network analysis by varying the alpha parameter.¹ If those counties with more distant ties gain more centrality power when the radius of assessment is increased, the type of social capital these counties possess would decrease crime.

As expected, higher levels of bridging social capital have crime mitigating characteristics. The steps in the procedure to measure long-range bridging social capital involve evaluating alpha centrality at a high absolute value of alpha ($\alpha = -0.1$), then evaluating alpha centrality in the same network at a low absolute level of alpha ($\alpha = -0.01$). In the first case centrality is evaluated at a high radius from the reference county. In other words, friend counties of friend counties (and their friend counties etc.) add to the reference county's centrality. In the second case, the radius is reduced. Only nearby counties, with social connections to the reference county, add to its' centrality. Lastly, The alpha centrality measures obtained from the low radius analysis (low absolute value of alpha) are subtracted from the alpha centrality measures obtained with the higher absolute level of alpha.

The resulting values of the difference in centrality measures contain the county-level gain in long-range bridging social capital from increasing the centrality assessment radius. The counties that gain the most in this exercise are the ones that contain the most long-range bridging social capital. Those counties that do not have much long range social connections do not gain centrality when the radius of assessment is increased. The gain in centrality from increasing the radius of centrality assessment (i.e. increasing the alpha parameter) mitigates crime levels. Substituting the gain-in-centrality values for the combined social capital variable in the crime regressions shows that long range bridging social capital reduces

¹Ghosh et al. (2011) and Bonacich (1987) describe this beneficial feature of alpha centrality. Bonacich uses the letter β for the attenuation parameter instead of α .

county-level crime rates.

This paper is organized as follows. In the next section I briefly describe the literature on social capital. First I list some important studies in the development of the concept of social capital, especially as they pertain to the importance of accounting for both bonding and bridging social capital. Also, I will describe several works on the measurement of social capital. The method section will describe the way alpha centrality is used to arrive at a comprehensive measure of social capital. In the data section I will describe the data sets used in the empirical validation of the new social capital measure. The results section will describe my empirical findings, followed by the conclusion.

2.2 Descriptions and Literature

2.2.1 Definitions of Bonding and Bridging Social Capital

Social capital has two components: bonding and bridging social capital. While both types consist of social links between individuals, bonding social capital consists of links between people who share at least some friends in common. Granovetter (1973) noticed that the overlap in mutual friends between two connected individuals had a high correlation to the type and quality of link these two had. People with many friends in common had a qualitatively stronger connection, characterized by high levels of trust. Bridging social capital, on the other hand, is made up of links between two individuals who share no friends in common. These links connect two people who don't know each other well. While these acquaintances do not have a close relationship, built on high levels of trust, Granovetter (1973) noticed these links are nevertheless important to bring in new information into well-bonded groups, facilitating consensus creation and safeguarding against insular communities.

In the historical development of the concept of social capital the focus has been on what constitutes social capital. One of the first users of the term, Loury (1976), notes that social capital is the link between one's social position, and the benefit or disadvantage it

may confer, on the ability to acquire human capital. Humans are not born into a social vacuum. The type and level of social connections they get embedded into is responsible for a number of outcomes that concern economists: income, income inequality, and wealth acquisition. According to Loury, social connections act as a moderating factor in the level of human capital acquisition. A disciplined study of economic outcomes requires one to account for the social network an individual is embedded in. The recent availability of social connectivity data from the likes of Facebook and Twitter lets one take a new look at this old question.

Most researchers agree that social capital is an attribute of groups of people and it is found in the connections between people. However, according to Coleman (1988) it is desirable for well-connected networks to possess as many connections as possible.² According to him, well-bonded groups of people can build high levels of trust, safeguard against malfeasance, and promote cooperative behavior. While most agree on the benefits of close-knit, well-connected networks, many researchers (Granovetter (1973), Putnam et al. (2001)) maintain that high levels of group cohesion, or bonding social capital, are only desirable up to a certain point. Beyond this point, Adler and Kwon (2002) claim that too much bonding social capital can lead to parochialism and inertia. In an entrepreneurial environment, Uzzi (1997) notes that too much social capital can stifle innovation. Bobo (1988) notes that people in closely connected groups without contact or knowledge about other groups of divergent racial, ethnic or class backgrounds will develop and maintain prejudices against these other groups.

Weak ties are the building blocks of the other component of social capital — bridging social capital. Granovetter (1973) defines weak ties as connections between people who have no friends in common. While two connected people in a closely bonded group may share most of their friends in common, a pair of acquaintances may have no common friends. The tie that bonds these acquaintances is a weak tie. Weak ties, or bridging social capital, are instrumental in consensus formation, efficient information flow, and being able to access

²From a network perspective this corresponds to the strong ties described by Granovetter (1973).

resources in general. Burt (1992), calling the missing links around weak links structural holes, maintains that people familiar with these holes can exploit them for entrepreneurial gain. Weak ties bring in relevant information from other groups. Communities with too much bonding and not enough bridging social capital tend to be insular. This is the reason why it is important to incorporate bridging social capital into any measure of social capital.

2.2.2 Current measurement of bridging and bonding social capital

There have been attempts at differentiating between these two different types of social capital in empirical research. Beyerlein and Hipp (2005) set out to demarcate bonding from bridging capital by using the prevalence of different religious traditions in a certain geographic unit. Their argument is that different religious organizations get involved with outside groups at varying degrees. Beugelsdijk and Smulders (2009, 2003) and Sabatini (2005) count links between family and close friends as bonding social capital and the prevalence of women's sports and religious organizations as evidence of bridging social capital. While family and close friends do create bonding social capital, associations have both bonding and bridging social capital. Not only might different associations have various degrees of both types of social capital, but these levels might vary across different locations.

There have been several other studies with ad hoc characterizations of which type of social association might contain higher levels of bridging social capital. Some studies, like Hotchkiss et al. (2022) divide associations into those that contain higher levels of bridging social capital and those that have higher levels of bonding social capital based on the clubs' activities and focus. These ad hoc ways of classifying local social associations miss a key feature of bridging social capital. It is the weak links that connect divergent groups. Different social associations will have varying degrees of out-of-group connections; however, it is ineffective to simply classify some as containing more bridging social capital. Some types of clubs will have a high degree of heterogeneity in their level of bridging social capital based on what area of the country or the county they are in. For example, the socioeconomic makeup of a

chess club will depend on whether it is in the South of the US or the West. This in turn will play a role in the amount of bridging social capital these clubs create.

Coffé and Geys (2007) measure the composition of demographic traits in social associations to discern various levels of bridging versus bonding social capital. They calculate the prevalence of a certain demographic category and measure how much this deviates from the prevalence of the trait in the locale where the association meets. The more representative an association is of key demographic traits of the community, the more bridging capital it has. If a certain trait is overrepresented, the group is thought of as having too many people of a certain type. This group might not be as inclusive of others and is thought to have a lower level of bridging social capital. For example, if a club is made up of 75% of 65 year olds or older but the community where the club is located only has 20% of this demographic, this club would be classified as low bridging capital. Age and sex are some variables along which the authors measure demographic deviation from community averages.

Hoyman et al. (2016), using the method developed by Coffé and Geys (2007) to differentiate between bonding and bridging social capital, find a positive relationship between bridging social capital and per capita income levels. Their county-level analysis finds that, as the number of bridging organizations goes up by 1%, per capita income goes up by 0.014%. They find no statistically significant relationship between bridging social capital and income inequality. Furthermore, the authors find no correlation between bonding social capital and either the level of per-capita income or income inequality.

My study is more disciplined in measuring the potential of the existence of bridging capital of different social organizations. However, it is still only measuring potential levels of bridging social capital. Friends can be of various ages and of either sex. There are definitely several generations of people in families as well as different sexes. These compositions can be replicated in social clubs as well. The nature of the club might determine it's age composition: sports clubs might have more young people, while some other clubs might have more people from a certain sex. This approach can pick up the potential of the presence of bridging

links. The method would have to be refined to more accurately deduce the presence of bridging social capital. The over- or under-representation of certain demographic groups in a club is not necessarily indicative of the presence or lack of diversity. The divergence of these demographic traits from community averages would have to be conditional on the traditional make-up of the clubs being examined. The most accurate measurement of bridging social capital would involve the counting of links between otherwise closed-off social groups.

In a departure from Coffe and Geys, Kyne and Aldrich (2020) measure bridging social capital by the rate of membership (per 10,000 persons) in religious, civic, fraternal and union organizations. They use race and ethnic fractionalization, income, educational inequality, and gender inequality as proxies for bonding social capital. In their work, a place with complete racial or ethnic dissimilarity (as measured on a scale from 0 to 1) would contribute to low bonding capital. While Coffee and Geys use club-level deviations from community means to measure bridging social capital, Kyne and Aldrich use a similar measure to measure bonding social capital. Bonding social capital is by definition local as it is measuring local clustering. Bridging social capital consists of the links between these clusters. Since these two measures are inherently different, they have to be measured differently — ideally from different sources.

2.2.3 Latest developments in network-derived social connectivity studies

The network-centrality-based method that allows the combination of differently measured variables into one comprehensive social capital measure is alpha centrality. The details of this endeavor are described in the method section. To glean an idea of its usefulness it is helpful to look at some of the previous works that have used alpha centrality or network topology to analyze social connections.

Sajuria et al. (2015) analyze the formation of bonding and bridging social capital from Twitter data by looking at the Occupy Wall street movement, the Chilean Presidential Election and the Enough Food For Everyone IF campaign in London, England. They use

the Burt (1992) concept of structural holes to analyze network constraint and the formation of bridging social capital, in addition to using Coleman’s closure concept to look at the formation of bonding social capital. Sajuria et al. (2015) use network formation concepts to construct hypothetical networks randomly,³ with preferential attachment⁴ and finally using elements of both randomness and preferential attachment.⁵ The networks of connections obtained from the Twitter data then are compared to the artificially created networks. Sajuria and co-authors find that, while there is organic formation of bonding social capital in all three social movements, bridging social capital only forms in those where there are organizers involved, such as the Chilean Presidential Election and the IF Campaign.

Overbey et al. (2013) also use data from Twitter. Using Twitter data they analyze social influence of individual Twitter users. More specifically they look at the tweets connected to the 2011 Egyptian Revolution. Using re-tweets to construct a social lattice (in other words an adjacency matrix of directed social connections) and using the number of people who saw an individual user’s message as a measure of influence exogenous to the social network, they are able to pinpoint the most influential people, via Twitter, of the 2011 Egyptian Revolution. Alpha centrality allows such an analysis since it extracts an eigenvector of network influence in addition to incorporating an exogenous measure — in this case the number of people who have seen messages of the Twitter users whose level of influence is being evaluated. While this work is not looking at social capital it looks at a related measure — social influence and consensus creation. In the present paper individual-level data of social capital is not available. Alpha centrality can be used to get around this data insufficiency — distant county-to-county social connections can be used to impute bridging social capital while within-county associational membership can be used as an external-to-network measure to count bonding social capital.

³Erdős and Rényi (1959) detail the formation of random networks based on a certain probability of link formation.

⁴Nodes connect to each other with a probability based on how many links they each have as in Barabási and Albert (1999).

⁵In Watts and Strogatz (1998), nodes form randomly into several circular structures and then each node has its links probabilistically rewired to other nodes.

2.3 Data

To develop the new measure of social capital that contains both bonding and bridging social capital, I use several datasets. For bonding social capital, I use the dataset of Rupasingha et al. (2006). This dataset contains county-level measures of association memberships, voter turnout and the response rates to the US Census. The authors also extract a principal component from these three variables to arrive at a county-level social capital measure.

For bridging social capital, I use county-to-county social connections data from Facebook. Facebook released the Social Connectedness Index in 2016. It is a snapshot of normalized county-to-county Facebook links. The dataset contains every US county’s social links aggregated to the county level. This allows me to treat this data set as an adjacency matrix.⁶ Therefore, I can extract an eigenvector from it. The elements of this vector serve as rankings of network centrality of the individual counties. These rankings constitute the relative level of bridging capital across all the counties.

I combine bridging capital with bonding capital through alpha centrality as described in the next section. Alpha centrality is a more general way of calculating eigenvector centrality. The combined measure serves as my principal variable of interest. Since the bonding social capital measure is released for 1990, 1997, 2005, 2009 and 2014, I am able to construct the combined social capital measure for these years.⁷

To utilize the new social capital variable I first estimate hate-group formation for 2005 using the same model as Goetz et al. (2012) have done for 2008. I also utilize data on various social and economic outcomes: hate-group formation, various crime levels, obesity and unemployment rate. Hate-group data comes from the Southern Poverty Law Center’s Hate Group Map. Crime-level data is from The Federal Bureau Of Investigations Uniform Crime Reporting database. The unemployment rate is from NIH IPUMS. I obtained obesity

⁶More specifically, it is a social graph.

⁷The SCI is only released for the 2016 Facebook connections. Bailey et al. (2018b) claim that Facebook connections reveal real life social connections. As such, it is reasonable to assume the SCI links are valid for times when Facebook did not exist. This assumption implies that bridging social capital is stable over time.

rates from the Federal Communication Commission Broadband Health dataset.

Data on election outcomes comes from the MIT Election Lab. The data on the adherents of Mainline, Evangelical and Catholic religions come from American Religious Data Archives. The source of Median Income Percent Black, Percent Hispanic Percent over 65 old, Median Property Tax, Foreign Born Rate and Unemployment Rate is NIH IPUMS. The MSA dummy was constructed from Bureau of Labor Statistics data. The divorce rate variable is from Bowling Green University's County-Level Marriage & Divorce Data, 2010.

The data on hate-group formation is from 2005. The outcome variable is the number of hate-groups across counties. This data comes from the Southern Poverty Law Center's Hate Group Map. It takes on discrete values. In the Goetz et al. (2012) study there are five basic categories of covariates that are thought to affect the presence of hate-groups in a county.⁸ (1) The covariates that are included in the history and geography category are population and state fixed effects⁹ (2) Some hate-groups are formed as a response to perceived government presence. Median Property Tax is used to proxy for government interference. (3) Frustration is measured by educational attainment and the unemployment rate¹⁰. (4) The proportion of the population that is Black and the proportion of the population that is Hispanic and Foreign Born Rate to proxy for status anxiety. (5) Divorce rate is used to measure social disintegration.

Furthermore, as in the Goetz et al. (2012) study, I add variables that measure the rate of adherence to three major religious groups: Mainline Protestant, Catholic, and Evangelical denominations. Income levels are included because, according to Goetz, people may join a hate-group if they feel there is competition for their earning potential, or because counties with higher levels of incomes will have more resources to alleviate hate-group formation. Finally the number of Walmart stores in a county is included to proxy for alienation. Goetz

⁸Goetz et al. (2012) indicate that the variables in these five categories were deemed as determinants of hate-group presence by Jefferson and Pryor (1999).

⁹Goetz et al. (2012) use a dummy variable for whether the a county is located in a state that was part of the Confederacy.

¹⁰In the present study, unemployment was flagged by a Variance Inflation Factor test as having excessive multicollinearity.

et al. (2012) hypothesize that Walmarts, by driving out small town businesses, reduced the number and wealth of the local middle class. This in turn might lead to alienation and increased propensity to form or join a hate-group.

2.4 Method

The nature of this study requires the discussion of two methods. Since the innovation of this paper rests in the introduction of a new type of social capital measure, it is imperative to start with the description of how this new measure is constructed. I detail this process in 4.1 - 4.3. In 4.4 I detail the estimation of marginal effects of the social capital variable on various outcome variables using OLS.

2.4.1 Measuring Social Capital

Traditionally, social capital measures are tallied up from responses to survey questions or counted from memberships in social organizations at some geographic level. The survey questions typically try to establish the level of trust of the person being interviewed. The surveys ask questions about how trustworthy the interviewee finds other people, other racial or ethnic groups, or the government that their community is embedded in. One of the weaknesses of these surveys is that they can get expensive as the sample size increases. The other one is that people tend to answer questions to impress the interviewer. Glaeser et al. (2000) notes this propensity where people rated themselves as more generous in a survey than was actually measured in an experiment.

The second method of measuring social capital for a certain community is to count the number of memberships in local social organizations. These organizations include sport clubs, religious organizations, and civic and educational groups. Since social capital is meant to measure the social links between people, membership in these organizations is a good proxy for social capital. People voluntarily join in these associations to either socialize with each

other, solve social challenges, or seek spiritual enlightenment in an organized manner. In all these settings members establish social links with each other. Therefore membership data in these social groups serves as a good proxy for bonding social capital.

Bonding social capital is the term used for connections between people characterized by strong ties. The defining characteristics of these social ties is that people tend to know each other well, and they share many mutual friends among each other. For example, if person A is socially connected to person B and the set S of A 's friendship network has a significant or perfect overlap with set T of B 's friendship network, the two people have a strong tie between them. These ties make up bonding social capital. The principal source of bonding social capital is the family unit. In communities where people seek out each other's company, the level of bonding social capital is higher. In these communities, people can achieve a consensus on certain issues more quickly.

As described in the literature section, in recent years some studies have emerged pointing out possible negative consequences of too much bonding social capital.¹¹ When people cluster into tight-knit groups they may have the propensity to exclude others, fall for harmful ideologies, and have more difficulties achieving consensus within a wider geographic setting. Bridging social capital may thus be key to keeping communities connected to each other.

Using network terminology, it is important to have links between well-connected clusters of nodes or components. While bonding social capital between two individuals is characterized by an overlap in their set of social connections, bridging social capital is defined by no overlap of sets of friends between two linked individuals. If friends A and B from above do not have any common friends and there is no overlap in their friendship sets S and T , the tie between them is a weak link. Weak links are the building blocks of bridging social capital. In common terms, A and B would be considered acquaintances. According to Granovetter (1973), the strength of the tie between A and B would become stronger as they start to develop mutual friends. Partly, this is because, as their friendship sets develop more overlap,

¹¹Satyanath et al. (2017) finds that municipalities with higher social club memberships succumbed to the National Socialist party propaganda at a faster rate in pre-WWII Germany.

the time the two friends spend around each other increases probabilistically. At some point, tension would develop if A and B did not become close friends.

While friendships are important for maintaining healthy communities, so are acquaintances. People can obtain novel information through weak links. Well-connected groups might have the same information bouncing around within the group.

2.4.2 Using Social Network Data to Measure Bridging Social Capital

One obstacle in measuring bridging social capital stems from the nature of it. While it is relatively easy to obtain membership data from social organizations, it is more difficult to measure weak links, or find proxies for social links between acquaintances. Survey-based research also has trouble uncovering weak links between the interviewed since by nature these are people who are not well known to the subject. This is where the Social Connectedness Index from Facebook becomes helpful. The SCI, previously described in the data section, provides the normalized number of friendship links between US counties. While the within-county friendship links contain a large number of strong ties such as the links between close friends and family, the number of links between counties contains a large number of weak links. Therefore the county-to-county SCI links can be used to construct a measure of bridging social capital.

Before moving on to the construction of the combined variable, it is useful to analyze why the more geographically distant ties may contain more weak ties. Geographically distant ties occur between people who have moved a long time ago or recently. If the move occurred recently, the network ties between the mover and her old network might stay the same; however, the mover will establish new links in the new locality. These new ties will be weak ties for the mover initially. Once the mover has met more people in the new locality, she will be more likely to share friendship circles with new people she meets in her new place. However, from the perspective of the people from where this person moved, the new people she meets will constitute weak ties. Geographically distant ties will contain more weak ties

than proximate ties.

Figure 2.1 demonstrates the instantaneous (at a moment in time) logic behind why distant social connections contain more weak ties. The nodes designated by letter a through x are individuals. In close proximity, a will have friends b and c who are also friends between themselves. As one looks at more distant radii of the social connections of a , links between friends of a have to decrease. As Bailey et al. (2018b) observe, the elasticity of friendship to distance is negative. This structure necessitates the higher likelihood of weak links as one considers links between successively more distant individuals. The only way more and more distant friends of a could also be friends of each other is for a 's elasticity of friendship to distance to be positive. While there might be individual instances where this is the case, overall the elasticity is -1.99 .¹²

Iyer et al. (2005) note that mobility has a negative association with social capital. People who do not expect to live in a place for a long time will not make the necessary investment to build their social network. In this paper, I am looking at distant ties and not mobility. It is reasonable to believe that there might have been people moving between two distant counties that have a high level of social ties between them; however, these movers can still expect to live in the place they had moved to for a long time. As such, they still will make the necessary time commitments to build up their social network. In addition, the social capital Iyer et al. (2005) refer to when they note the negative association between mobility and social network investments is bonding social capital. The availability of county-to-county social links from Facebook makes the detection of weak links possible on a country-wide scale.

¹²According to Bailey et al. (2018b) the elasticity of friendship to distance is -1.99 within 200 miles and -1.16 over distances of 200 miles.

2.4.3 *Alpha centrality - the mechanics of combining the two types of social capital*

Ideally, a measure of social capital should contain the relevant levels of both types when measured in a certain locality, such as a county. The membership levels in local associations has been used as a proxy by researchers from times even before Putnam et al. (2000)'s work on social capital. More recently, there have been attempts at differentiating between the two types of social capital.¹³ The measurement of the two types of social capital depends on ad-hoc decisions by researchers.

The county-to-county social connections can be used to detect the levels of bridging social capital across counties. Since the SCI dataset contains every county's social connections, finding a network centrality measure for each county would serve as a measure of bridging social capital. Eigenvector centrality is an appropriate method for this purpose. According to Bonacich (1987) eigenvector centrality uses the concept that a connection to a node which itself has higher connectivity contributes to the centrality of the focal node more than a connection to a node with fewer connections. To get around the simultaneous determination of centrality, an eigenvector is extracted from the adjacency matrix of social connections. The values in the eigenvector provide a network centrality ranking for the members of the network.

$$A^T x = \lambda x \tag{2.1}$$

Where A is the adjacency matrix from the social graph, x is the eigenvector and λ is the largest eigenvalue.

Since the SCI provides a weighted adjacency matrix for US counties, the eigenvector

¹³Coffé and Geys (2007) measure the composition of demographic traits in social associations to discern various levels of bridging versus bonding social capital. Kyne and Aldrich (2020) measure bridging social capital by the rate of membership (per 10,000 persons) in religious, civic, fraternal and union organizations and use race and ethnic fractionalization, income, educational, and gender inequality as proxies for bonding social capital. Beugelsdijk and Smulders (2009) count links between family and close friends as bonding social capital and the prevalence of women's sports and religious organizations as evidence of bridging social capital.

provides a ranking of the network centrality of each county. If, as I hypothesized above, the county-to-county connections contain a higher level of bridging social capital, then the eigenvector will rank counties by their level of bridging social capital. However, bridging social capital is only one side of the story. Bonding social capital is an integral component of total social capital. Alpha centrality is an enhanced form of eigenvector centrality.¹⁴ It allows the inclusion of a vector of exogenous variables into the eigenvector centrality formula. When this vector contains a row of ones, the result is eigenvector centrality. In the current study, I include the county-level bonding social capital values in the exogenous vector. Alpha centrality therefore will produce a county-level ranking based on both bonding and bridging social capital:

$$x = \alpha A^T x + e \tag{2.2}$$

Where x is the eigenvector as above, e is the exogenous vector of characteristics, and α is the attenuation factor, or how much network-derived centrality adds to status (centrality).

2.4.4 Measuring the impact of social capital in various outcomes

With the new, combined, social capital measure at hand, I estimate several outcomes using OLS regression. Several types of crime outcomes, unemployment and obesity are estimated using the combined social capital measure and several other relevant determinants of the relevant outcomes. For hate-group formation I use Poisson regression since the dependent variable, hate-groups, is a discrete number for each county. All of these specifications are compared to their counterparts using the predominant, bonding social capital measure used in the literature. As will become apparent, the new type of social capital measure is significant and is of the expected sign in more models than the bonding one. The new measure of social capital is an improvement over the old one because it combines both types of social capital and it measures them appropriately, based on their characteristics.

¹⁴Bonacich (2007) notes alpha centrality is similar to the Katz (1953) prestige measure.

2.5 Results

In this section, I present the results of several empirical models using the new, combined social capital variable. First, I present eigenvector centrality and alpha centrality values for the US state of Oklahoma. This exercise demonstrates the centrality ranking of Oklahoma counties based on their social connectedness, as provided by Facebook SCI data, and alpha centrality based on social connectedness and within county bonding social capital. Next, I present some intuitive, descriptive statistics on the different values of alpha. In this table I present county-level per-capita income, unemployment rate, and Democrat% - Republican% vote margin for the top and bottom deciles by alpha centrality rank for instances when alpha ranges from -0.1 to 0.1 in increments of 0.01. Following this, I estimate the hate-group formation model of Goetz et al. (2012) for the year 2005. I then estimate a number of county-level criminal outcomes, obesity and unemployment, with my preferred measure of social capital and the bonding social capital measure of Goetz et al. (2012) for 2016. Finally, I examine the impact of centrality gain of a county as an explanation for some counter-intuitive results in the previous subsections.

2.5.1 Descriptives and how alpha correlates with outcomes

Table 2.1 presents centrality values of Oklahoma counties. Eigenvector- and Alpha-centrality are used to rank Oklahoma's 77 counties. Eigenvector centrality only uses the county-to-county social connections that were provided by Facebook. Alpha centrality uses the social connections between counties and also within-county associational membership (bonding social capital). The eigenvector rankings are displayed in column *e_rank*. Number one in the ranking is Dewey county — a county in Northwest Oklahoma about one hundred miles from the capital, Oklahoma City. Ranked last is Adair County, bordering Arkansas. These rankings are based on social connections alone. It is worth noting that the out-of-state social connections of these 77 counties are not included in this analysis. This exercise is for

demonstrative purposes.

Alpha centrality incorporates within-county social capital measures as well. This ranking picks Tulsa, a metropolitan area of about one million people, as number one, Cleveland County, the home of the University of Oklahoma, as number two, and Oklahoma City, the capital of Oklahoma as number three. These rankings are more in line with expected levels. This exercise is merely an illustration of how network centrality measures can be used to rank counties based on social connectivity data.

Table ?? presents top and bottom deciles of per capita income, unemployment rate and vote margin for social capital measures of a few, selected levels of the centrality attenuation parameter alpha. For negative values of alpha, the counties in the top 10% by social capital have higher income levels, lower unemployment rates and they vote more Republican than counties in the bottom decile. When alpha is zero, this pattern in the values in the top and bottom deciles is more pronounced. At positive alpha's, counties in the upper deciles have lower incomes and higher unemployment. The vote pattern is the same as with negative alpha's.¹⁵

2.5.2 Hate-group formation estimation using combined social capital

Goetz et al. (2012) look at the effect of Walmart store and social capital on county-level hate-group formation. They build on Jefferson and Pryor (1999) by expanding the number variables that are thought to be responsible for the number of hate-groups in a certain county. Jefferson and Pryor (1999) included historical and geographical determinants of hate-group formation proxied by population density and whether a county belonged to the Confederacy.¹⁶ Government interference is measured by property taxes. Culture and possible discontent is measured by the education level. The proportion of the county population that is Black was used to measure status anxiety. Finally, they included the divorce rate to

¹⁵It is worthwhile to note that according to Bonacich (2007), positive alpha parameters evaluate a network as an information or social network (non-conservative, information propagation). Negative alpha evaluates a network as a conservative one, where a certain value or good is being subdivided and negotiated over.

¹⁶In the present study I use state fixed effects instead of the Confederate dummy.

proxy for social disintegration, and rural/urban status is used as a control. Goetz et al. (2012) add the number of Walmarts, social capital and the level of adherents to three main religious groups to these variables. The spread of Walmart into rural areas precipitated the decline of many businesses there, with the concomitant decline of the small-town middle class. Therefore it is thought that the number of Walmarts would facilitate hate-group formation. Social capital is thought to act as a social lubricant and reduce the number of hate-groups in a county. The effect of the level of religious adherence varies by the type of religion. Goetz et al. (2012) note that while Mainline and Catholic groups exercise more community outreach, Evangelicals are more inward-looking, hence they might facilitate hate-group formation while the first two groups would not.

The primary purpose of replicating Goetz et al. (2012) for a different year is to study my newly created combined social capital variable. Their study uses variables from 2008. As Goetz et al. (2012) do, I use Poisson regression to estimate the effect of the above listed variables on hate-group formation for 2005. Of most concern in the model is the effect of social capital. I will enumerate the effect of other variables if they diverge from the Goetz study of 2008 or the replicated model using their social capital variable.

When I replicate the Goetz et al. (2012) model with 2005 data, the effect of their social capital variable is insignificant as depicted in column 2 of Table 2.3. The new, combined social capital variable derived using alpha centrality is significant at the 1% level and is of the expected, negative, sign (column 1). The effect of Walmarts is similar; 0.01 in the regression with the combined social capital variable and 0.013 with the bonding social capital. The coefficient on religious groups is similar between the two years but differs from the Goetz et al. (2012) model using 2008 data. In 2005, both the coefficient on Catholic and Evangelicals is positive and significant at the 10% level. Goetz et al. (2012) hypothesize that since Catholics and Mainline Protestants exercise more social outreach, these religious groups inhibit the formation of hate-groups in their community. Evangelicals, however, are more inward looking and might not inhibit hate-groups. The coefficient on the county-level rate of Evangelical

adherence is positive and significant in their paper (using 2008 data), while the coefficient on Mainline is negative. The Catholic coefficient is insignificant in their estimation.

The coefficient on population is positive in all specifications; however, it is insignificant in the replicated estimation of 2005 data with bonding social capital and in the original Goetz estimation. The population coefficient in the model using the combined social capital measure (2005 data) is 0.28% and significant at the 1% level. Similarly, the Metropolitan Statistical Area coefficient, a dummy variable indicating an urban county, is also positive (0.92%, significant at 1%). This is similar in all the estimations; urban counties are more likely to host a hate-group than rural ones, holding fixed other variables.

The effect of the percentage of foreign born people in the county is also positive, as is the level of property taxes. All estimations concur on this. The effect of county-level divorce rate is positive in the Goetz paper; however, it is insignificant with 2005 data, whether the combined social capital variable is used or the bonding one.

The combined social capital variable seems to be measuring social capital with less noise than the bonding one. It is worth noting that a grid search algorithm was used to select an alpha that minimized the mean squared errors in the regressions. The algorithm picked out -0.01 as the appropriate value. As mentioned above, negative values of alpha indicate a conservative social network. This is a network where some set amount of value is percolating over the network and individuals are negotiating over it.¹⁷

2.5.3 Social Capital Effects On Crime

Next I examine the effect of social capital on crime. The outcome variables are total crime rate, violent crime, property crime and murder. As before, the new, combined, social capital measure is compared to the existing, bonding social capital measure as a determinant of crime. Since the determinants of crime are similar to the ones affecting hate-group formation,

¹⁷Positive alphas evaluate an information, or social network. In general alpha centrality is a steady-state solution to some type of diffusion, making it a metric for representing the process of information sharing or influence or division of a good over time.

the control variables remain the same as in the hate-group formation model (Walmart data was not available for 2016. Fixed effects are at the state level). Odd numbered columns of table 2.4 display the coefficients in the models where the alpha-centrality-derived, combined social capital measure is used. In all crime models, with the exception of property crime, the coefficients are positive (significant at the 1% level for total crime and violent crime; at the 5% level for murder). Property crime is not statistically significantly different from zero. The coefficient on bonding social capital does not attain significance in any of the models.

There are some notable differences between the crime and hate-group models. While the median property tax coefficient is positive (significance at the 1% level) in the hate-group model it is negative (significance at the 1% level) in all the crime models, combined and bonding social capital models alike. Jefferson and Pryor (1999) and Goetz et al. (2012) measured government interference using property taxes in their models. They explained the positive coefficient on property taxes by claiming more government interference propagated the formation of hate-groups. In the crime models, property taxes potentially measure the wealth level of the county; hence, higher taxes decrease crime. The coefficient on median income is also negative (significant at the 1% level for property crime and at the 5% level for murder). It is also possible that high-property-tax counties provide more public goods to their residents, which ameliorates their propensity for criminal acts.

Another point of divergence between the crime and hate-group models is the behavior of the foreign born rate variable. In the crime models, this coefficient is not significant in any type of crime model.¹⁸ In the hate-group model the coefficient on foreign born rate is positive at 1.24 (significant at the 1% level), the combined social capital model and 1.23 (significant at the 1% level) in the model with bonding social capital.

¹⁸A meta-study of 51 studies from 1994 – 2014 by Ousey and Kubrin (2018) concluded that the relationship between immigration and crime is negative but weak. The relationships that the 51 studies found were dependent on the research design.

2.5.4 *Social Capital Effects On Other Outcomes*

Table 2.5 displays the results for the county-level obesity and unemployment regressions. The grid search in these estimations picked out a zero level of alpha as the one minimizing mean squared errors. At $\alpha = 0$ the social capital variable is the same as the one supplied by Goetz et al. (2012); the one measuring county-level bonding social capital. For the unemployment model, the coefficient on social capital is negative and significant at the 1% level. For the regression estimating obesity, the coefficient on social capital is also negative and significant at the 5% level. In these models, incorporating bridging social capital increases mean squared errors in the estimating regressions.

2.5.5 *Centrality Gain From Increasing Network Radius*

In this section, I evaluate a county's gain in centrality from increasing the alpha parameter. The counterintuitive results from the crime models require further analysis. The coefficient on social capital for most types of crime variables is positive and significant in the crime models. This indicates that social capital increases crime. However, it is possible that the low network diameter — over which the low alpha parameter is assessing network centrality — is the culprit. The level of alpha was selected using a grid search that looks for a level of alpha that minimizes mean squared errors in the estimating equations. However, this level of alpha might be too small to evaluate county-level bridging capital. Recall equations 2.1 and 2.2 that introduced alpha centrality. Equation 2.2 can be transformed into Equation 2.4. Equation 2.4 can be expressed as an infinite summation in equation 2.5. The $\alpha^k * A^{k+1}$ terms are assessing centrality over successively higher paths. As the path length is increased α^k becomes attenuated as the exponent k grows in value. Centrality from higher path lengths adds less to the centrality of the reference node. This is the manner in which varying alpha can modify the diameter of centrality assessment. For values of alpha close to the maximum of $(1/\lambda)$, alpha centrality approaches eigenvalue centrality — centrality is assessed over the maximum number of path lengths. Decreasing alpha to zero reduces

alpha centrality to degree centrality as only the neighbors of the reference node add to its centrality.

$$x = \alpha A^T x + e \tag{2.3}$$

$$x = (I - \alpha A^T)^{-1} \times e \tag{2.4}$$

$$x = (A + \alpha A^2 + \alpha^2 A^3 + \alpha^3 A^4 + \dots) \times e \tag{2.5}$$

In the above equations as the length of paths is increasing the weight these nodes add to the centrality of the focal node is decreasing due to the attenuating parameter alpha and its exponent. $1/\alpha$ is the effective network radius that centrality is assessed over. Therefore those nodes that gain the most centrality when the level of alpha is increased have the greatest number of connections, or bridging social capital, at higher radii.

In this section I develop a variable that looks at which counties gain the most centrality when the level of alpha is increased from -0.01 to -0.1 . As the diameter of centrality assessment is increased, those counties that have more distant connections will gain centrality. As mentioned above, it is these counties that have individuals with high levels of bridging social capital. The alpha-centrality-derived social capital measure used above in the crime models might be looking at connections between groups in neighboring counties¹⁹. Communities where there are high levels of short-range bringing ties might foster criminal activity. The model that assesses the gain in centrality when alpha is increased tests for long-range bridging social capital.

As evident in Table 2.6, the more long-range bridging ties a county has, the lower the crime levels. The social capital variable of interest is composed of the difference between combined social capital at alpha= -0.1 and the measure with alpha= -0.01 . This variable

¹⁹Ballester et al. (2010) note that it is the removal of low-level criminals who connect criminal groups that has the highest possibility of disrupting criminal networks.

measures the contribution of bridging capital at the level of the global network minus the contribution of bridging social capital at the local level.²⁰

2.6 Conclusion

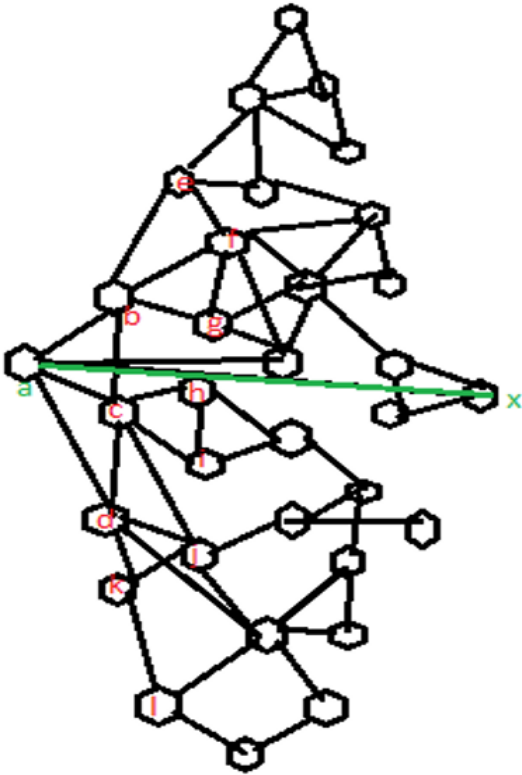
This paper examines the measurement of social capital. It looks at the importance of measuring both bonding and bridging components of it. Since these two components have inherently different structures, measuring them from geographical-area associational membership levels introduces the possibility of mis-measurement problems.

The paper then proposes alpha centrality to better measure social capital. Alpha centrality combines a network-derived centrality measure, which in this case proxies for weak ties that constitute bringing capital, and the within-county bonding social capital, measured from within-county associational membership levels. The new, combined variable is compared to bonding social capital in several models of crime, obesity and unemployment. Combined social capital is also compared to bonding social capital in a replication of the Goetz et al. (2012) hate-group formation study for 2005. Here also the updated combined measure performs better.

In a robustness check, the counterintuitive positive sign of the combined social capital variable is investigated by looking at the gain in county-level network centrality as the radius of network analysis is increased via the the alpha parameter. The resulting negative crime coefficients in the regressions point to the mitigating effects of long-range bridging social capital on all crime types.

²⁰When alpha is close to zero (with e a vector of ones) alpha centrality is reduced to the measure of degree centrality. With alpha approaching $1/\lambda$, alpha centrality measures eigenvector centrality.

Figure 2.1: Relationship between geographically distant ties and weak links



A's friends B, C and D are friends among themselves. At the next radius, E and F are also friends with A and among themselves. However, at this distance (radius) F and I are not friends with each other. As one looks at further radii from A there are fewer and fewer mutual friendships. As one goes farther, the chance for mutual friends diminishes. A and Y have no friends in common (weak tie). For A to keep having strong ties at further and further radii there would have to be more connectivity between friends as distance increases. This is not the case in general — elasticity of friendship to distance is decreasing.

Table 2.1: Eigenvector and alpha centrality status analysis of county level data of Oklahoma

fips	County name	Eigenvector centrality	Alpha centrality	e_rank	a_rank	fips	County name	Eigenvector centrality	Alpha centrality	e_rank	a_rank
40001	Adair	0.012316	0.000320	77	21	40079	Le Flore	0.026285	0.000857	65	9
40003	Alfalfa	0.551283	-0.000051	9	46	40081	Lincoln	0.043484	-0.001639	51	74
40005	Atoka	0.071083	-0.000045	34	44	40083	Logan	0.058570	-0.002126	43	76
40007	Beaver	0.296295	-0.000383	14	60	40085	Love	0.052960	0.000314	45	23
40009	Beckham	0.485219	-0.000331	11	58	40087	McClain	0.080576	0.000563	28	11
40011	Blaine	0.610668	0.000223	7	32	40089	McCurtain	0.034601	-0.000146	57	50
40013	Bryan	0.061965	-0.000277	42	55	40091	McIntosh	0.028817	-0.001567	62	73
40015	Caddo	0.139826	-0.000580	23	66	40093	Major	0.669800	0.000095	6	39
40017	Canadian	0.113324	0.000271	26	29	40095	Marshall	0.050435	0.000318	47	22
40019	Carter	0.071248	0.000141	33	35	40097	Mayes	0.022804	-0.001496	69	72
40021	Cherokee	0.025044	0.000300	67	26	40099	Murray	0.066833	-0.000503	39	63
40023	Choctaw	0.051715	-0.000055	46	47	40101	Muskogee	0.025211	-0.000125	66	49
40025	Cimarron	0.162171	0.000415	21	15	40103	Noble	0.097957	0.000325	27	20
40027	Cleveland	0.047839	0.002822	49	2	40105	Nowata	0.020817	0.000253	72	30
40029	Coal	0.079330	0.000108	29	38	40107	Okfuskee	0.030123	-0.000283	61	56
40031	Comanche	0.048246	0.000156	48	34	40109	Oklahoma	0.056027	0.002300	44	3
40033	Cotton	0.077245	0.000336	31	18	40111	Okmulgee	0.033627	0.000302	58	24
40035	Craig	0.030733	0.000042	60	42	40113	Osage	0.024315	-0.002305	68	77
40037	Creek	0.026304	0.001121	64	5	40115	Ottawa	0.022403	0.001038	70	7
40039	Custer	0.494240	0.000332	10	19	40117	Pawnee	0.128389	-0.000898	55	69
40041	Delaware	0.020448	0.000135	73	37	40119	Payne	0.128389	0.000571	24	10
40043	Dewey	1.000000	0.000135	1	36	40121	Pittsburg	0.043143	-0.000422	52	61
40045	Ellis	0.814160	0.000168	3	33	40123	Pontotoc	0.070319	0.000408	35	16
40047	Garfield	0.256513	-0.000548	16	65	40125	Pottawatomie	0.041967	-0.000736	53	68
40049	Garvin	0.062918	-0.000527	41	64	40127	Pushmataha	0.068957	0.000064	38	40
40051	Grady	0.068969	-0.001353	37	71	40129	Roger Mills	0.565565	-0.000253	8	54
40053	Grant	0.261638	0.000279	15	28	40131	Rogers	0.018346	-0.001969	74	75
40055	Greer	0.251545	-0.000371	17	59	40133	Seminole	0.041639	0.000367	54	17
40057	Harmon	0.180283	-0.000432	20	62	40135	Sequoyah	0.016665	0.000534	76	13
40059	Harper	0.688761	-0.000223	5	53	40137	Stephens	0.063357	-0.000672	40	67
40061	Haskell	0.032694	0.000045	59	41	40139	Texas	0.148723	-0.000291	22	57
40063	Hughes	0.047522	0.000548	50	12	40141	Tillman	0.127501	0.001062	25	6
40065	Jackson	0.199516	-0.000193	19	52	40143	Tulsa	0.022394	0.003107	71	1
40067	Jefferson	0.073175	-0.000172	32	51	40145	Wagoner	0.017402	0.001764	75	4
40069	Johnston	0.069032	-0.000047	36	45	40147	Washington	0.026604	0.000283	63	27
40071	Kay	0.078615	0.000938	30	8	40149	Washita	0.344577	0.000237	13	31
40073	Kingfisher	0.357426	-0.000952	12	70	40151	Woods	0.715263	-0.000106	4	48
40075	Kiowa	0.243966	0.000437	18	14	40153	Woodward	0.875939	0.000018	2	43
40077	Latimer	0.037570	0.000300	56	25						

¹ Observations at the county level. Only connections for Oklahoma are used to construct the adjacency matrix. The social connections of Oklahoma counties with out-of-state counties are ignored in this exercise. Social connectivity data was obtained from Social Connectedness Index released by Facebook Corporation. e_r , a_n and a_r , a_n are the centrality rankings based on eigenvector centrality an alpha centrality respectively.

Table 2.2: Simple correlations between alpha centrality and outcomes for different values of attenuation parameter alpha

	alpha = -.1		alpha = -.05		alpha = -.01		alpha = 0	
alpha centrality	lower 10%	upper 10%	lower 10%	upper 10%	lower 10%	upper 10%	lower 10%	upper 10%
per capita income (10,000 dollars)	2.79	2.83	2.81	2.82	2.78	2.76	2.09	2.85
unemployment rate	0.0662	0.0654	0.0667	0.0655	0.0676	0.0676	0.0912	0.0381
vote margin	-0.1914	-0.2004	-0.1842	-0.2052	-0.1887	-0.2162	-0.3183	-0.4723
	alpha = .01		alpha = .05		alpha = .1			
alpha centrality	lower 10%	upper 10%	lower 10%	upper 10%	lower 10%	upper 10%		
per capita income (10,000 dollars)	2.81	2.77	2.83	2.78	2.82	2.78		
unemployment rate	0.0641	0.0661	0.0647	0.0663	0.0653	0.0661		
vote margin	-0.1925	-0.1976	-0.1903	-0.2013	-0.1933	-0.1983		

¹ Per capita income, Unemployment rate, vote margin of the top and bottom decile of alpha-centrality derived social capital. Alpha parameter varies in increments of 0.01 from -0.1 to 0.1. Observations at the county level. Social connectivity data was obtained from Social Connectedness Index released by Facebook Corporation. Vote margins are calculated from county level % of vote for Democrat minus the % for Republican candidate. This data comes from the MIT Election Lab. The source of per capita income and unemployment rate is NIH IPUMS.

Table 2.3: Effect of Social Capital on Hate Group Formation

	Hate	
	(1) Combined (-0.01)	(2) Bonding only
Social capital	-0.0487*** (0.0172)	-0.00317 (0.0873)
Number of Walmarts	0.0989*** (0.0186)	0.131*** (0.0148)
Catholic rate	0.000822* (0.000482)	0.000875* (0.000487)
Mainline rate	-0.000525 (0.000524)	-0.000637 (0.000563)
Evangelical rate	0.000880* (0.000508)	0.000890* (0.000521)
Percent Black	7.165*** (0.977)	7.520*** (0.972)
Percent Black squared	-8.712*** (1.671)	-9.089*** (1.674)
Population (million)	0.282*** (0.103)	0.0445 (0.0571)
Percent High school or more	7.702*** (1.087)	7.350*** (1.133)
Divorce rate	0.000891 (0.00350)	0.000901 (0.00356)
MSA	0.920*** (0.143)	0.914*** (0.144)
Foreign born rate	1.240*** (0.401)	1.229*** (0.402)
Median property tax (dollar)	0.000402*** (0.0000821)	0.000424*** (0.0000822)
Median income (dollar)	-0.0000235*** (0.00000660)	-0.0000235*** (0.00000662)
Constant	-9.597*** (0.954)	-9.405*** (1.037)
N	3035	3035

Notes: Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01. ¹ Combined Social Capital is calculated using alpha-centrality. The social capital values are presented in units of standard deviation and are demeaned. The network component of it was calculated using Social connectivity data from Social Connectedness Index released by Facebook Corporation. The county-level social capital to calculate combined social capital comes from Goetz et al. (2012). The source of Median Income Percent Black, Percent Hispanic Percent over 65 old, Median Property Tax, Foreign Born Rate and Unemployment Rate is NIH IPUMS. The MSA dummy was constructed from Bureau Of Labor Statistics data. Catholic, Mainline and Evangelical rates are from Association of Religious Data Archives. The Divorce Rate variable is from Bowling Green University's County-Level Marriage & Divorce Data, 2010. State Fixed Effects are used in all models. State Fixed Effects subsume the Confederate dummy used by Goetz et al. (2012). Since the dependent variable is count data, in both models a Poisson Regression Model is used.

Table 2.4: Effect of Social Capital on Crime

	Total crime			Violent crime			Property crime			Murder		
	(1) Combined (-0.01)	(2) Bonding	(3) Combined (-0.01)	(4) Bonding	(5) Combined (-0.05)	(6) Bonding	(7) Combined (-0.01)	(8) Bonding	(9) Combined (-0.01)	(10) Bonding	(11) Combined (-0.01)	(12) Bonding
Social capital	58.23*** (15.85)	-5.789 (26.46)	51.84*** (6.527)	8.440 (10.99)	9.615 (10.88)	-14.22 (18.41)	0.328** (0.154)	0.203 (0.257)	0.328** (0.154)	-14.22 (18.41)	0.328** (0.154)	0.203 (0.257)
Percent Black	641.3*** (177.7)	624.5*** (180.7)	186.4** (73.17)	155.2** (75.05)	455.6*** (123.9)	469.2*** (125.8)	14.74*** (1.730)	14.37*** (1.758)	14.74*** (1.730)	469.2*** (125.8)	14.74*** (1.730)	14.37*** (1.758)
Percent Hispanic	122.6 (204.8)	108.9 (205.3)	73.47 (84.35)	64.50 (85.27)	48.65 (142.9)	44.33 (142.9)	3.301* (1.995)	3.280 (1.997)	3.301* (1.995)	44.33 (142.9)	3.301* (1.995)	3.280 (1.997)
Percent over 65 years old	-1356.2*** (479.2)	-1308.1** (514.2)	-164.4 (197.4)	-215.9 (213.5)	-1190.5*** (334.3)	-1092.3*** (357.9)	-8.746* (4.667)	-10.11** (5.000)	-8.746* (4.667)	-1092.3*** (357.9)	-8.746* (4.667)	-10.11** (5.000)
Percent with a bachelor degree over 25 years of age	1475.3*** (294.6)	1469.6*** (302.3)	-48.25 (121.3)	-86.88 (125.5)	1518.8*** (205.5)	1556.4*** (210.4)	13.26*** (2.869)	12.65*** (2.940)	13.26*** (2.869)	1556.4*** (210.4)	13.26*** (2.869)	12.65*** (2.940)
Population (million)	5225.7*** (57.95)	5279.9*** (56.20)	1973.9*** (23.87)	2023.8*** (23.34)	3251.2*** (39.77)	3256.2*** (39.12)	35.14*** (0.564)	35.47*** (0.547)	35.14*** (0.564)	3256.2*** (39.12)	35.14*** (0.564)	35.47*** (0.547)
Median property tax (dollar)	-0.304*** (0.0381)	-0.315*** (0.0384)	-0.136*** (0.0157)	-0.143*** (0.0160)	-0.167*** (0.0266)	-0.172*** (0.0268)	-0.00233*** (0.000371)	-0.00234*** (0.000374)	-0.00233*** (0.000371)	-0.172*** (0.0268)	-0.00233*** (0.000371)	-0.00234*** (0.000374)
Divorce rate	-0.138 (1.333)	-0.133 (1.337)	0.0317 (0.549)	0.0163 (0.555)	-0.168 (0.930)	-0.150 (0.930)	0.00659 (0.0130)	0.00628 (0.0130)	0.00659 (0.0130)	-0.150 (0.930)	0.00659 (0.0130)	0.00628 (0.0130)
MSA	130.4*** (38.84)	126.7*** (39.34)	22.23 (16.00)	21.93 (16.34)	108.3*** (27.10)	104.8*** (27.38)	0.757** (0.378)	0.787** (0.383)	0.757** (0.378)	104.8*** (27.38)	0.757** (0.378)	0.787** (0.383)
Foreign born rate	53.10 (120.5)	50.75 (120.8)	-2.174 (49.65)	-3.184 (50.18)	55.11 (84.09)	53.94 (84.10)	0.0823 (1.174)	0.0878 (1.175)	0.0823 (1.174)	53.94 (84.10)	0.0823 (1.174)	0.0878 (1.175)
Catholic rate in 2010	-0.121 (0.171)	-0.114 (0.171)	-0.124* (0.0702)	-0.123* (0.0712)	0.00210 (0.119)	0.00901 (0.119)	-0.00163 (0.00166)	-0.00170 (0.00167)	-0.00163 (0.00166)	0.00901 (0.119)	-0.00163 (0.00166)	-0.00170 (0.00167)
Mainline rate	0.00166 (0.132)	0.0114 (0.140)	0.0249 (0.0543)	0.00945 (0.0582)	-0.0226 (0.0919)	0.00198 (0.0976)	0.000925 (0.00128)	0.000562 (0.00136)	0.000925 (0.00128)	0.00198 (0.0976)	0.000925 (0.00128)	0.000562 (0.00136)
Evangelical rate	0.249* (0.147)	0.250* (0.151)	0.0487 (0.0607)	0.0332 (0.0628)	0.201* (0.103)	0.217** (0.105)	-0.0000345 (0.00143)	-0.000317 (0.00147)	-0.0000345 (0.00143)	0.217** (0.105)	-0.0000345 (0.00143)	-0.000317 (0.00147)
Median income (dollar)	-0.00312 (0.00246)	-0.00280 (0.00248)	0.0000337 (0.00101)	0.000180 (0.00103)	-0.00317* (0.00171)	-0.00298* (0.00172)	-0.0000533** (0.0000239)	-0.0000539** (0.0000241)	-0.0000533** (0.0000239)	-0.00298* (0.00172)	-0.0000533** (0.0000239)	-0.0000539** (0.0000241)
Constant	-132.3 (221.0)	-147.3 (238.0)	-28.59 (91.03)	2.839 (98.84)	-103.4 (154.2)	-150.0 (165.7)	0.538 (2.153)	1.231 (2.315)	0.538 (2.153)	-150.0 (165.7)	0.538 (2.153)	1.231 (2.315)
N	3019	3019	3019	3019	3019	3019	3019	3019	3019	3019	3019	3019

¹ Combined Social Capital is calculated using alpha-centrality. The social capital values are presented in units of standard deviation and are de-meaned. The network component of it was calculated using Social connectivity data from Social Connectedness Index released by Facebook Corporation. The county-level social capital to calculate combined social capital comes from Goetz et al. (2012). The source of Median Income Percent Black, Percent Hispanic Percent over 65 old, Median Property Tax, Foreign Born Rate and Unemployment Rate is NIH IPUMS. The MSA dummy was constructed from Bureau Of Labor Statistics data. Catholic, Mainline and Evangelical rates are from Association of Religious Data Archives. The Divorce Rate variable is from Bowling Green University's County-Level Marriage & Divorce Data, 2010. State Fixed Effects are used in all models.

Notes: Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 2.5: Effect of Social Capital on Obesity and Unemployment

	Unemployment rate $\times 100$		Obesity rate	
	(1) Combined (0)	(2) Bonding	(3) Combined (0)	(4) Bonding
Social Capital	-0.475*** (0.0645)	-0.475*** (0.0645)	-0.140** (0.0702)	-0.140** (0.0702)
Percent Black	6.237*** (0.441)	6.237*** (0.441)	9.724*** (0.480)	9.724*** (0.480)
Percent Hispanic	-0.134 (0.501)	-0.134 (0.501)	-3.806*** (0.545)	-3.806*** (0.545)
Percent over 65 years old	-0.823 (1.254)	-0.823 (1.254)	-15.72*** (1.364)	-15.72*** (1.364)
Percent with a bachelor degree over 25 years of age	-3.639*** (0.737)	-3.639*** (0.737)	-22.53*** (0.802)	-22.53*** (0.802)
Population (millions)	-0.106 (0.137)	-0.106 (0.137)	-0.575*** (0.149)	-0.575*** (0.149)
Median Property tax (100 dollars)	0.412*** (0.0937)	0.412*** (0.0937)	-0.374*** (0.102)	-0.374*** (0.102)
Divorce Rate (per 100 people)	-0.602* (0.326)	-0.602* (0.326)	-0.232 (0.355)	-0.232 (0.355)
MSA	-0.134 (0.0959)	-0.134 (0.0959)	0.105 (0.104)	0.105 (0.104)
Foreign Born Rate	-0.252 (0.295)	-0.252 (0.295)	-0.219 (0.321)	-0.219 (0.321)
Catholic Rate (per 100 people)	-0.0815* (0.0418)	-0.0815* (0.0418)	-0.0273 (0.0455)	-0.0273 (0.0455)
Mainline Rate (per 100 people)	-0.234*** (0.0342)	-0.234*** (0.0342)	0.0894** (0.0372)	0.0894** (0.0372)
Evangelical Rate (per 100 people)	-0.174*** (0.0369)	-0.174*** (0.0369)	-0.0318 (0.0401)	-0.0318 (0.0401)
Median Income (in \$100s)	-0.1000*** (0.00604)	-0.1000*** (0.00604)	0.00312 (0.00657)	0.00312 (0.00657)
Constant	13.90*** (0.580)	13.90*** (0.580)	39.23*** (0.632)	39.23*** (0.632)
N	3019	3019	3019	3019

¹ Combined Social Capital is calculated using alpha-centrality. The social capital values are presented in units of standard deviation and are de-meanned. The network component of it was calculated using Social connectivity data from Social Connectedness Index released by Facebook Corporation. The county-level social capital to calculate combined social capital comes from Goetz et al. (2012). The source of Median Income Percent Black, Percent Hispanic Percent over 65 old, Median Property Tax, Foreign Born Rate and Unemployment Rate is NIH IPUMS. The MSA dummy was constructed from Bureau Of Labor Statistics data. Catholic, Mainline and Evangelical rates are from Association of Religious Data Archives. The Divorce Rate variable is from Bowling Green University's County-Level Marriage & Divorce Data, 2010. State Fixed Effects are used in all models. Notes: Standard errors are reported in parentheses. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Table 2.6: Radius Analysis by varying α

	(1)	(2)	(3)	(4)	(5)	(6)
	Total Crime	Violent Crime	Property Crime	Murder	Unemployment	Obesity
Δ social capital as $\alpha \uparrow$ 0to-0.1	-58.28*** (15.86)	-51.89*** (6.532)	-6.380 (11.06)	-0.328** (0.154)	-0.0000571 (0.000391)	-0.0409 (0.0422)
Percent Black	641.3*** (177.7)	186.4** (73.17)	454.8*** (123.9)	14.74*** (1.730)	0.0567*** (0.00438)	9.573*** (0.473)
Percent Hispanic	122.6 (204.8)	73.48 (84.35)	49.09 (142.9)	3.301* (1.995)	-0.000196 (0.00505)	-3.764*** (0.545)
Percent 65 or Older	-1356.2*** (479.2)	-164.4 (197.4)	-1191.8*** (334.3)	-8.746* (4.667)	-0.0412*** (0.0118)	-16.70*** (1.275)
Percent 25 or Older with a Bachelor Degree Population/Million	1475.4*** (294.6)	-48.22 (121.3)	1523.5*** (205.5)	13.26*** (2.869)	-0.0481*** (0.00726)	-22.86*** (0.784)
Median Property Taxes	5225.7*** (57.95)	1973.9*** (23.87)	3251.8*** (40.43)	35.14*** (0.564)	-0.000568 (0.00143)	-0.597*** (0.154)
Divorce Rate 2010	-0.304*** (0.0381)	-0.136*** (0.0157)	-0.168*** (0.0266)	-0.00233*** (0.000371)	0.0000508*** (0.00000939)	-0.000340*** (0.000101)
MSA	-0.138 (1.333)	0.0317 (0.549)	-0.169 (0.930)	0.00659 (0.0130)	-0.0000668** (0.0000329)	-0.00252 (0.00355)
Foreign Born Rate	130.4*** (38.84)	22.22 (16.00)	108.1*** (27.10)	0.757** (0.378)	-0.000310 (0.000958)	0.137 (0.103)
Catholic Rate	53.10 (120.5)	-2.173 (49.65)	55.28 (84.10)	0.0823 (1.174)	-0.00214 (0.00297)	-0.207 (0.321)
Mainline Rate	-0.121 (0.171)	-0.124* (0.0702)	0.00215 (0.119)	-0.00163 (0.00166)	-0.0000103** (0.00000421)	-0.000338 (0.000454)
Evangelical Rate	0.00166 (0.132)	0.0249 (0.0543)	-0.0232 (0.0919)	0.000925 (0.00128)	-0.0000319*** (0.00000325)	0.000646* (0.000351)
Median Income 2016 (10,000 dollars)	0.249* (0.147)	0.0487 (0.0607)	0.200* (0.103)	-0.0000345 (0.00143)	-0.0000232*** (0.00000363)	-0.000486 (0.000392)
Constant	-312 (24.6)	0.337 (10.1)	-31.5* (17.1)	-0.533** (0.239)	-0.0105*** (0.000606)	0.0151 (0.0654)
N	-132.3 (221.0)	-28.59 (91.03)	-103.6 (154.2)	0.538 (2.153)	0.155*** (0.00545)	39.69*** (0.588)
	3019	3019	3019	3019	3019	3019

Notes: Standard errors are reported in parentheses. * p <0.1, ** p <0.05, *** p <0.01.

¹ Combined Social Capital is calculated using alpha-centrality. The social capital values are presented in units of standard deviation and are de-meaned. The variable of interest is the difference between social capital calculated with zero alpha and alpha = -0.1. The network component of it was calculated using Social connectivity data from Social Connectedness Index released by Facebook Corporation. The county-level social capital to calculate combined social capital comes from Goetz et al. (2012). The source of Median Income Percent Black, Percent Hispanic Percent over 65 old, Median Property Tax, Foreign Born Rate and Unemployment Rate is NIH IPUMS. The MSA dummy was constructed from Bureau Of Labor Statistics data. Catholic, Mainline and Evangelical rates are from Association of Religious Data Archives. The Divorce Rate variable is from Bowling Green University's County-Level Marriage & Divorce Data, 2010. State Fixed Effects are used in all models. State Fixed Effects subsume the Confederate dummy used by Goetz et al. (2012).

Chapter 3

Heterogeneity in Child

Quantity-Quality Trade-off: A

Machine Learning Approach

3.1 Introduction

Economists have long been interested in understanding the relationship between family size and child quality. This is not only because family environment is a primary component to child's quality (Black et al. 2005), but also because understanding this relationship is important to policy makers. The theoretical quantity-quality model (Becker and Lewis 1973, Becker and Tomes 1976) predicts a negative effect of family size on child quality. This quantity-quality trade-off has become the main justification for family planning campaigns to increase population quality by curbing population growth. We use recently developed machine learning methods to investigate the heterogeneous effects of family size. We find the mother's age at first birth, her income and education levels play a large role in the negative effect of additional children on health outcomes.

Many empirical studies have tested the quantity-quality trade-off using data from various countries. They find mixed results. Some studies find that child quality is not significantly affected by family size (Kessler 1991; Guo and VanWey 1999; Black et al. 2005; Angrist et al. 2010; Zhong 2014). Other studies find evidence of a quantity-quality trade-off (Rosenzweig and Wolpin 1980; Cáceres-Delpiano 2006; Li et al. 2008; Liu 2014; Rosenzweig and Zhang 2009). A third set of studies suggests a positive effect of family size on child quality (Lee 2008, Qian 2009).

Testing for the existence of a quantity-quality trade-off is complicated by the endogeneity of family size. If parents who place a greater value on child quality more prefer fewer children, then the relationship between family size and child quality will be driven by parental preferences rather than family size. A commonly used method to establish the causal effect of family size on quality is to use instrumental variables (IV). Empirical studies have used the birth of twins (Rosenzweig and Wolpin 1980; Black et al. 2005; Cáceres-Delpiano 2006; Li et al. 2008; Angrist et al. 2010), twinning by birth order (Rosenzweig and Zhang 2009), child sex composition (Angrist et al. 1998; Angrist et al. 2010; Lee 2008), and the One Child Policy (OCP) (Qian 2009; Liu 2014) to extract exogenous variation in family size. One of the distinguishing features of the IV approach is that, without a homogeneity assumption, the IV estimate is the local average treatment effect (LATE) of a group of compliers (Angrist and Imbens 1995). It is possible that the compliers of different instrument variables are so disparate that the average treatment effects of compliers using different instrument variables differ from each other significantly. The mixed results from using the IV approach suggests that there might be heterogeneity in the effect of family size on child quality.

There are several theories purporting to explain the potential heterogeneity in the impact of family size on child quality. Rosenzweig and Wolpin (1980) point out that whether there is a quantity-quality trade-off depends on whether child quantity and child quality are substitutes or complements. If parents are facing a budget constraint on how much they can invest in total child quality, then having an additional sibling reduces the average amount

of resources invested in each child, hence reducing individual child quality. However, it is possible that, for some families, additional siblings benefit from existing children by stabilizing the parental relationship (Becker 1998, Black et al. 2005), increasing the likelihood of the mother staying at home to provide child care (Gelbach 2002; Black et al. 2005; Ruhm 2008), or generating other positive spillovers (Bandura and Walters 1977). It is also possible that some parents adjust to an exogenous increase in family size by working longer hours, consuming less leisure, or investing less in themselves rather than decreasing quality inputs on each child (Angrist et al. 2010).

While the relationship between family size and child quality has been under intensive investigation in the literature, few studies have looked beyond LATE. This paper uses recently developed machine learning methods to examine the potentially heterogeneous effect of family size on child quantity in the Chinese context. We find substantial heterogeneity in treatment effects of additional children on child quality. The heterogeneity in treatment effects is most pronounced for girls. Giving birth later in a mother’s reproductive life cycle, being more educated and having higher incomes all contribute to treatment effect heterogeneity.

This paper makes several contributions to the literature. First, we move beyond LATE and allow for the heterogeneous treatment effects of family size on child quality. Moreover, as measures of quality, we focus not only on children’s education but also on children’s health. Examining the heterogeneity of treatment effects of family size on child quality is crucial to understanding the process of human capital investment within households. It also has strong policy implications as family planning policies might have different impacts on households with disparate characteristics.¹

Second, the machine learning method we propose advances the traditional approach in several ways. Traditional approaches to explore heterogeneous effects involve analyzing subgroups and including interaction terms in the model. These methods require the researcher

¹For related studies please consult Hedrich (2011) and Brinch et al. (2017).

to have a thorough understanding of the research question to define subgroups. Additionally, the potential for cherry-picking problems may arise and unexpected subgroups may be missed (Lee et al. 2020). This approach involves interacting the treatment with various covariates, usually one covariate at time. This raises the probability of spurious findings (Davis and Heller 2017). Commonly used linear approaches may also fail to capture some nonlinear treatment effects since there might be nonlinearities in child quality from changes in family size (Løken et al. 2012). In contrast, machine learning methods use data-driven algorithms to detect heterogeneity which avoids leaving out important heterogeneities. Furthermore, flexible modeling alleviates the concern of the existence of nonlinear relationship between family size and child outcomes. Therefore this method makes the discovery of treatment heterogeneity more reliable for policy determination. We find that being a girl, mothers age at first birth, parental education and income levels play a larger role in the quantity-quality trade-off. Hence, policy-makers can fine-tune what segment of the population they need to concentrate on to alleviate the detrimental effects of having more children on child quality.

Third, this paper joins the literature on the effects of the OCP on child outcomes. It fills an important gap that exists in the literature by investigating the heterogeneity the OCP's treatment effects. As such we use data from the 1976 Statistical Yearbook of China where we collect information on provincial characteristics prior to the implementation of the One-Child-Policy (OCP). We use this to control for fertility preference prior to the OCP. We also use data from the China Health and Nutrition Survey (CHNS), a longitudinal dataset that collects detailed information on households on health, education, and a variety of other demographic information.

The paper is structured as follows: the next section will discuss the data used in our estimation. Following this, we lay out the empirical methodology in section 3.3. In section 3.4 we explain the results, including the IV and the estimation results of machine learning algorithm. Section 3.5 concludes.

3.2 Data

The first source of data we use is from the 1976 Statistical Yearbook of China where we collect information on provincial characteristics prior to the implementation of the One-Child-Policy.

The second data source we use in this paper is from the China Health and Nutrition Survey (CHNS), a longitudinal dataset that collects detailed information of households on health, education, and other basic demographic information. From this dataset we calculate height-for-age z-score (HAZ) and weight-for-age z-score (WAZ) as the two measures for child’s health. HAZ and WAZ are constructed using the British 1990 Growth Reference. For child’s education, we use relative educational attainment, and a dummy variable for being currently enrolled in school.²

In our analysis of the quantity-quality trade-off, we restrict our sample to first-born children. We impose a few restrictions to our sample. First, we exclude children from households with twins because the birth of twins results in shorter birth spacing. This could potentially interfere with our main results. Second, we restrict our sample to children aged between 6 and 17 at the time of survey. Before age 6, children are not in school and after age 17, the family influence is diminishing. Third, we restrict our sample to children who were born in 1976 or after. Before the adoption of the OCP, the family planning policy was known as the “later(marriage), longer(intervals), fewer(children)” (Chen and Huang 2018). The recommended birth spacing is 4 years (Liu 2014). Since the OCP was first carried out in some provinces in 1979, for couples who followed the previous family planning policy and had their firstborn in 1976 or after, their decision to have a second child would be heavily affected by the implementation of the strict OCP. Finally, we exclude firstborns whose province of residence is not available. In addition, we drop the observations from Chongqing because we were unable to obtain the pre-OCP provincial characteristics of Chongqing as Chongqing

²Following Rosenzweig and Wolpin (1980), relative educational attainment is constructed as $Educ_{igw}/\overline{Educ}_{gw}$, where $Educ_{igw}$ is the years of schooling of child i at age g from wave w and \overline{Educ}_{gw} is the average years of schooling of children at age g in wave w .

was a part of Sichuan province at that time.

We present summary statistics in Table 3.2. The firstborns in our sample are relatively evenly distributed between boys and girls. 90% of the children are currently enrolled in school. The average age of the sample is 11 years and the average years of schooling is 5.48, implying that, on average, children in our sample are in their last year of primary school. We break down our sample into boys and girls. We highlight three findings from comparing boys' and girls' summary statistics. First, boys and girls are very similar in terms of family background, such as parental health and parental educational attainment. Second, we do not observe large disparities between boys' and girls' education measures. Third, HAZ and WAZ diverge between boys and girls. Girls' HAZ and WAZ are much lower.

3.3 Empirical Methodology

3.3.1 Heterogeneity in exposure to the One Child Policy

The One Child Policy was introduced to curb rapid population growth in China. Although it was a national policy, its implementation varied across provinces. First, the year of implementation of the OCP differed. The official document about the implementation of the OCP was released in 1979 but the actual implementation year among provinces ranged from 1979 to 1984. Second, provinces made modifications to the OCP over the years. The timing and the content of the modifications are often different. For example, the OCP faced strong resistance among rural households, particularly among those whose firstborn was a girl (Zhang 2017). Later, many provinces relaxed their OCP to allow rural households to have a second child. Among these provinces, the year of the rule relaxation ranged from 1985 to 1998. In some provinces, all rural households were eligible for this relaxation, while in other provinces, only rural households with a female firstborn were permitted to have a second child. Several provinces, such as Jiangsu, Zhejiang, Jilin, never carried out this relaxation. Table 3.1 presents the details of heterogeneity in the OCP implementation.

In addition to the divergence in the OCP implementation, the birth year of the mother also contributes to heterogeneity in the household’s exposure to the OCP. The fertility of households with a mother whose prime fertility years (21 – 35) were not covered by the OCP is less likely to be influenced. If a mother entered her prime fertility years after the implementation of the OCP, then the longer she was exposed to the OCP, the greater impact the OCP would have on her fertility decision.

3.3.2 *Instrumental Variable*

Substantial heterogeneity in a household’s exposure to the OCP generates a unique source of exogenous variation in family size. Borrowing the idea from Huang (2021), we use the share of a mother’s prime fertility years covered by the strict OCP as the instrumental variable:

$$\text{coverage} = \frac{\text{number of prime fertility years covered by the strict OCP}}{\text{total number of prime fertility years}}$$

The strict OCP is defined as only one child being allowed per couple. Suppose province A started implementing the OCP in 1979. Starting from 1985, eligible households were allowed to have a second child. Now we consider a household in province A with a mother who was born in 1954. The mother entered her prime fertility years in 1975 and exited in 1989. If the household was not eligible for having a second child, then 11 (from 1979 to 1989) of the mother’s most fertile years were exposed to the strict OCP. *coverage* for this household would be 11/15. If the household qualified for the OCP relaxation, then only 7 (from 1979 to 1985) years of the mother’s most fertile years were subject to the strict OCP. In this case, *coverage* would be 7/15. The larger *coverage* is, the more influence the OCP has on the household’s fertility choice. We should expect there to be a negative relationship between *coverage* and number of siblings.

3.3.3 *Validity of IV*

Our IV exploits the provincial variation in the OCP policies to account for exogenous variation in family size. One main concern is that variation in the OCP regulations is associated with provincial characteristics that may also influence fertility. For example, provinces that carried out the OCP relaxation policies sooner might have a stronger preference for a larger family. To alleviate this concern, we include pre-OCP provincial characteristic variables in the model to control for the potential correlation between the OCP policies and preexisting provincial fertility preferences.

Another concern is that couples might manipulate their eligibility for the OCP relaxation to have more children. The relaxed OCP mainly allowed two types of couples to have more than one child: ethnic minority couples and rural couples whose first child was a girl. An individual who desired more children could secure the eligibility by marrying a member of an ethnic minority. Past studies (Huang and Zhou 2015) show that the OCP induced more inter-ethnic marriages; however, it is less concerning in our analysis because nearly all of our sampled provinces required both spouses to be ethnic minority to be qualified for the OCP relaxation. The only exception is Guangxi province, which allowed couples with one ethnic minority spouse to have a second child. However, its policy only lasted for a short period of time before it was tightened to require both spouses to be an ethnic minority.

Rural couples could manipulate their eligibility status by utilizing sex-selective abortion to choose the gender of their first child. However, they had little incentive to do so. Rural couples who were pregnant with a boy were unlikely to terminate the pregnancy due to the strong son preference. Rural couples who were pregnant with a girl had little incentive to terminate their pregnancy because they were allowed to have a second child.

The validity of our IV strategy hinges on the assumption that family size is the only channel through which OCP policies affect child quality. Empirical studies (Ebenstein 2010, Chen et al. 2013) show that the OCP leads to sex ratio distortion among high-order births. This is concerning to our analysis because this means that the OCP might also affect sibling

sex composition. Sibling sex composition has its own effect on child quality as parents might allocate resources differently between boys and girls (Behrman et al. 1986).

To examine the empirical relationship between sibling sex composition and our IV, we estimate the following regression:

$$SibSex_{ipt} = \beta_2 coverage_i + D_{it}\gamma_3 + C_p\delta_3 + w_t + \eta_{ipt} \quad (3.1)$$

where $SibSex_{ipt}$ is a measure of sibling sex composition. In this analysis, we use two measures of sibling sex composition: an indicator for a male second birth and the fraction of male siblings. D_{it} is a set of individual characteristics, including child's age, mom's age at first birth, mom's age at first birth squared, parental education-level dummies. C_p is a set of pre-OCP provincial characteristic variables.

Consistent with Ebenstein (2010), we only find a significant relationship between sex composition and our IV for girls. Results in Table 3.3 show that, among households with female firstborns, mothers who are always ineligible for having a second child are 12.3% less likely to have a male second birth compared with those who are always eligible. Firstborns from an always ineligible household also have a lower fraction of male siblings. One possible explanation for the significant negative signs in column (3) and (4) is that ineligible couples have more incentive to use sex-selective abortion on their first birth, while eligible couples have more incentive to use it on their second birth. In the female first-born subsample, it is possible that the ineligible couples have a weaker son preference and less incentive to use sex selective abortion on their second birth than those eligible couples.

Results in Table 3.3 suggest that our IV estimates of the effect of family size on child quality for firstborn boys are not confounded by sibling sex composition. However, the IV estimates for firstborn girls should be interpreted with caution. It is plausible to assume that firstborn girls with a male sibling are allocated with less resources than those with a female sibling. As ineligibility of having a second child leads to lower chances of having a

male sibling, our IV estimates for firstborn girls are biased downward. We should consider those as the lower bound of the true estimates of the effect of family size for firstborn girls.

3.3.4 Local average treatment effect

We use a 2SLS approach to access the local average treatment effect of family size on quality of children. In the first stage, number of siblings is instrumented by *coverage* using equation (3.2).

$$nsib_{ipt} = \beta_1 coverage_i + D_{it}\gamma_1 + C_p\phi_1 + w_t + v_{ipt} \quad (3.2)$$

where $nsib_{ipt}$ is the number of siblings child i from province p in wave t . D_{it} contains a vector of individual control variables; C_p is set of provincial control variables³; w_t is wave fixed effects.

In the second stage, we use equation (3.3).

$$Y_{ipt} = \theta \widehat{nsib}_{ipt} + D_{it}\gamma_2 + C_p\phi_2 + w_t + \epsilon_{ipt} \quad (3.3)$$

where Y_{ipt} is the quality measure of child i in province p in wave t ; \widehat{nsib}_{ipt} is estimated from the first stage; the other variables are the same as defined in equation (3.2).

We consider child quality in terms of both health and education. For health, we use HAZ and WAZ. For education, we use relative educational attainment and a dummy variable for being currently enrolled in school.

The individual control variables include parental years of schooling, child's age, mother's age at first birth, and mother's age at first birth squared. In analysis focused on child's health, we also include parental height and parental weight. We replace the missing values of parental height/weight with the average sample parental height/weight. In addition, we include dummy variables for missing parental height/weight in case parents with some

³Note that we do not include province fixed effect because our IV is constructed based on provincial variations in the OCP implementation. Including province fixed effect would absorb some of the useful variation in the IV.

missing characteristics in the data set are characteristically different from those parents with no missing variables.

Provincial control variables include sex ratio, birth rate, log of GDP per capita, share of non-agricultural population, share of primary industry in GDP, and share of secondary industry in GDP to control preexisting provincial features such as fertility preferences.

3.3.5 *Conditional local average treatment effect*

The goal of this paper is to study heterogeneity in quantity-quality trade-offs. To do so, we estimate the conditional local average treatment of number of sibling on children’s quality using the instrumental forest within the generalized random forest framework developed by Athey et al. (2019). The method we use is an adaptive nearest neighbor approach, of which the weighting is derived from a random forest technique.

We are interested in estimating the following model for individual i , $i = 1, \dots, n$:

$$Y_i = \tau(X_i)W_i + \mu(X_i) + \epsilon_i \quad (3.4)$$

where Y_i is a child quality measure of individual i ; W_i is the treatment, in our analysis, it is the number of siblings individual i has; $\tau(X_i)$ captures the causal effect of the number of siblings on the quality of i ; $\mu(X_i)$ is a nuisance parameter; and ϵ_i is a noise term that is correlated with W_i . Because of the correlation between ϵ_i and W_i , we need to use an instrumental variable Z_i to generate a consistent estimate for $\tau(X_i)$. In this paper, we use *coverage _{i}* , the share of a mother’s prime fertility years covered by the strict OCP, as our instrumental variable. X_i contains the same covariates as in the 2SLS approach along with wave dummies.

$\tau(X_i)$ is obtained by minimizing equation (3.5), an empirical version of two moment functions based on the exclusion assumption of the instrumental variable: $\mathbb{E}[Z_i(Y_i - \tau(x)W_i - \mu(x))|X_i = x] = 0$ and $\mathbb{E}[Y_i - \tau(x)W_i - \mu(x)|X_i = x] = 0$.

$$(\hat{\tau}(x), \hat{\mu}(x)) \in \underset{\tau(x), \mu(x)}{\operatorname{argmin}} \left\{ \left\| \sum_{i=1}^n \alpha_i(x) \begin{pmatrix} Y_i - \tau(x)W_i - \mu(x) \\ Z_i(Y_i - \tau(x)W_i - \mu(x)) \end{pmatrix} \right\|_2 \right\} \quad (3.5)$$

The resulting instrumental forest estimator can be written as :

$$\hat{\tau}(x) = \frac{\sum_{i=1}^n \alpha_i(x)[Z_i - \bar{Z}(x)][Y_i - \bar{Y}(x)]}{\sum_{i=1}^n \alpha_i(x)[Z_i - \bar{Z}(x)][W_i - \bar{W}(x)]} \quad (3.6)$$

where $\alpha_i(x)$ is some kind of similarity weights that measure the relevance of individual i to the estimation of $\tau(x)$; $\bar{Z}(x) = \sum_{i=1}^n \alpha_i(x)Z_i$; $\bar{Y}(x) = \sum_{i=1}^n \alpha_i(x)Y_i$; $\bar{W}(x) = \sum_{i=1}^n \alpha_i(x)W_i$.⁴

Similarity weights, $\alpha_i(x)$, are derived from the random forest technique. In the traditional random forest method (Breiman (2001)), many trees are grown in the forest and each terminal leaf of the trees is associated with a specific prediction value. To obtain a prediction for a point of interest of x , x is pushed down through each tree till it hits a terminal leaf and a prediction for x from each tree is observed. The final prediction for x is done by averaging over predictions from all the trees in the forest. Each tree is grown by recursively splitting a random subset of covariate space and each split is chosen to maximize the prediction accuracy of the tree.

Instrumental forest, instead of looking for predictions, counts how often individual i ends up in the same terminal leaf with x among all trees. Similarity weights, $\alpha_i(x)$, are calculated as the following:

$$\alpha_i(x) = \frac{1}{B} \sum_{b=1}^B \frac{I\{X_i \in \mathcal{N}_b(x)\}}{|\mathcal{N}_b(x)|}$$

where B is the number of trees in the forest; $\mathcal{N}_b(x)$ is the number of individuals that fall into the same terminal leaf as x in tree b ; $I\{X_i \in \mathcal{N}_b(x)\}$ is an indicator function that

⁴To improve the performance of the instrumental forest in practice, Athey et al. (2019) suggest using the centered outcome \tilde{Y}_i , centered treatment \tilde{W}_i , centered instrumental variable \tilde{Z}_i to replace Y_i , W_i , Z_i , respectively. $\tilde{Y}_i = Y_i - \hat{y}^{(-i)}(X_i)$, where $\hat{y}^{(-i)}(X_i)$ is leave-one-out estimate of the marginal expectation of Y_i , computed without using the i th observation. \tilde{T}_i and \tilde{Z}_i are computed likewise. $\bar{Y}(x) = \sum_{i=1}^n \alpha_i(x)\tilde{Y}_i$, likewise for $\bar{W}(x)$ and $\bar{Z}(x)$.

takes on value 1 if individual i ends up in the same terminal leaf as x in tree b . The more frequently individual i falls into the same terminal leaf as x , the higher value $a_i(x)$ receives for individual i in the estimation of $\tau(x)$.

Compared to the traditional random forest, two other features of the instrumental forest used in this paper are worth highlighting. First, the splitting criterion is different. The traditional random forest focuses on delivering predictions; as a result, each split is chosen to maximize prediction accuracy. In the instrumental forest, each split seeks to maximize the heterogeneity in treatment effects across partitions. This splitting criterion helps improve the expected accuracy in predicting treatment effects (Athey et al. 2019). Second, trees grown in the instrumental forest are “honest” trees, meaning that a tree is constructed by using one subsample, while similarity weights derived from this tree are estimated using a different subsample. In other words, $\alpha_i(x)$ is obtained by using trees constructed without individual i . This method is similar in structure to cross-fitting.

We conduct our analysis in R, using the package `grf` developed by Athey et al. (2019). The instrumental forest estimator is obtained by the function `instrumental_forest`. Without formal criteria to guide our choices of the parameters in the function, we set all the parameters at their default values except two parameters: the number of trees and the minimum number of observations in each leaf. Increasing the number of trees reduces Monte Carlo error at the price of increasing computational cost. We choose to grow 15000 trees because the improvement of the estimate stability dramatically slows down when more trees are grown.⁵ The choice of the minimum number of observations in each leaf shows a trade-off between bias and variance. Large minimum size of each leaf produces overly simplified tree models that generate less heterogeneity. Previous studies typically set the minimum size between 1 and 10 (Davis and Heller 2017; O’Neill and Weeks 2018; Baiardi and Naghi 2020). In this paper, we choose 10 as our minimum leaf size.

⁵We use the median prediction variance as the measure for estimate stability. Figures 3.1 and 3.2 display the relationship between the number of trees and estimate stability.

3.4 Results

3.4.1 *Local average treatment effect*

Results from 2SLS estimations are reported in Table 3.4 and 3.5. First, we find that an additional sibling has a significantly negative impact on children’s WAZ. An increase in the number of siblings reduces first-born boy’s WAZ by 1.066 standard deviations.⁶ It reduces first-born girl’s WAZ by 0.816 standard deviations. The magnitude of the effect is smaller for girls. Second, we only find a significant effect of sibling size on HAZ among boys. One additional sibling decreases boy’s HAZ by 1.251. Third, all the coefficients on *coverage* are highly significant, meaning *coverage* is significantly related to the number of siblings. The F statistics are above the Stock-Yogo critical value of 10 % maximal IV size, which confirms that our IV is not weak.

Now we turn to the results of estimations using educational measures as dependent variables. We only find a significant impact of sibling size on relative educational attainment among female firstborns.⁷ An increase in sibling size leads to a 0.149 decrease in relative educational attainment for girls. We do not find any significant effect of additional siblings on the likelihood of school enrollment for girls or boys. The fact that sibling size has a limited effect on children’s education can be explained by the adoption of the Compulsory Education Law in 1986. The law requires that all children attend school for a minimum of nine years. Our results suggest that the Compulsory Education Law does protect children from the negative additional sibling effect on their education during their early school years.

3.4.2 *Heterogeneity in local average treatment effects*

The generalized random forest procedure adopted from Athey et al. (2019) allows for the estimation of heterogeneity in treatment effects. For girls, the distribution of treatment

⁶WAZ and HAZ are z-scores therefore the units of measurement are in standard deviations.

⁷Relative educational attainment is measured as a fraction of the mean group educational attainment. The group for calculating the mean is defined by age and wave.

effects has a relatively tight dispersion. There is relatively wider treatment heterogeneity for boys, as visible in Figure 3.3. The range of treatment effects of additional siblings is much more spread out for boys. This pattern is much more apparent for the distribution of the effects of an additional sibling on health variables than on educational variables.

Another interesting finding is that the distributions of treatment effects on HAZ and WAZ exhibit different patterns among boys and girls. For girls, the distribution of treatment effects on HAZ centers around zero, while almost the entire distribution of treatment effects on WAZ falls below zero. For boys, the distribution of treatment effects on HAZ is relatively similar to the distribution of treatment effects on WAZ. HAZ captures long-term health consequences, while WAZ presents short-term health consequences. This finding suggests that most of the girls experience a negative effect of an additional sibling in the short run, but some of them are able to catch up and even benefit from having more siblings in the long run.

In addition to the distributions of treatment effects, we also report quartile treatment effect means using all four dependent variables in Table 3.6. For boys, the quartile means range from -3.63 to 0.31 for HAZ, and from -2.40 to 0.95 for WAZ. For girls, the quartile means range from -1.26 to 0.71 for HAZ, and from -1.94 to -0.24 for WAZ. The ranges of quartile means for educational variables are much smaller. For boys, the quartile means range from -0.47 to 0.15 for relative education, and from -0.29 to 0.31 for school enrollment. For girls, the quartile means range from -0.83 to 0.29 for relative education, and from -0.21 to 0.03 for school enrollment.

The next step before examining treatment effect heterogeneity is to see whether the generalized random forest algorithm actually captures treatment heterogeneity in general. To perform this simple test we examine whether the subgroup that is hypothesized to have the largest treatment effect does in fact have a larger treatment effect than the other subgroup. In the spirit of Davis and Heller (2017), we first divide our sample into two subsamples based on the treatment effects estimated by the generalized random forest approach. One

subsample contains individuals whose treatment effects are above the median treatment effect. The other subsample contains the rest of the individuals. We then repeat the 2SLS regression discussed in Section 3.3.4 within each subgroup and compare the coefficients on sibling size across two subgroups. We present results in Table 3.7.

Consistent with the result in the graphs depicting the distribution of treatment effects we find the largest heterogeneity in treatment effects for the two health-related dependent variables — HAZ and WAZ. In the 2SLS regressions estimating HAZ and WAZ for the subgroup with the largest treatment effects additional siblings reduce health outcomes the most. Furthermore, the difference between the subgroup with the largest treatment effects and the smallest treatment effects are significant for both sexes and both dependent variables. In contrast, for the education outcome variables the treatment effect for the most negative treatment effect subgroup is insignificant as is the difference between this group’s treatment effect and the rest of the sample. The results of this simple check of the heterogeneity of treatment effects is in accordance with the graph in Figure 3.3. The results suggest that the generalized random forest algorithm is only able to capture the heterogeneity in treatment effects when the dependent variables are health related variables. Since we do not have enough evidence of whether the algorithm can also capture heterogeneity in treatment effect for education related variables, we focus our heterogeneity analysis on health related variables only.

In Table 3.8 we look at covariate means for the below-median and above-median treatment groups in estimating HAZ, separately for the two sexes. For both boys and girls, children who experience a larger negative effect of sibling size are from households with lower parental health measures, lower parental education, and mother giving birth at an older age. When the dependent variable is WAZ, the differences in coefficients display remarkable similarity, especially for the boys.

The generalized random forest algorithm provides statistics on variable importance. The number of times a variable is used to split the covariate space can be used to determine

its importance in the estimation of treatment heterogeneity. Athey et al. (2019) prove that splitting to minimize mean squared errors in the leaf nodes is equivalent to maximizing treatment heterogeneity.⁸

In Figure 3.4 the bars indicate the relative importance of variables for HAZ and WAZ estimation for girls and boys separately. As expected, father’s and mother’s height variables play an important role in growing the regression trees in the random forest algorithm. Furthermore, since height and weight are correlated both of these variables play an important role in estimations of both of the health-related child quality estimations — HAZ and WAZ. Parental education and mother’s age at first birth also play an important role in all estimations. This has important implications from a policy perspective. Governments can implement incentives to improve the education level of the population and help women to give birth during the most optimal times of their life-cycle.

Parental education

In Figure 3.5 and 3.6 we examine the heterogeneity of treatment effects looking at their distribution, color-coded by quartiles, for different levels of parent’s education. Some clear patterns emerge in these graphs. When HAZ is the dependent variable, increasing either parent’s educational level moves the distribution to the right — the quantity-quality trade-off is attenuated for higher-educated mothers and fathers of girls alike. There is no clear pattern for distributional shifts for boys. It is possible that as parental educational levels go up, there are more resources to invest in child quality.⁹ However, if parents set aside more resources to invest in boys regardless of income level, this channel is muted.

It is also of import that when the dependent variable is WAZ the pattern described above is not observed. The reason for this might be the difference between the two biometrics. Children’s weights respond faster to changes in nutrition, while height is slower to respond and it

⁸Athey et al. (2019) also use a penalty structure for variance for the treatment and control outcomes in the leaves.

⁹Chen and Li (2009) find positive effects of mother’s education on height-for-age z-scores using Chinese data

is thought to be a repository of long term nutritional investments. In other words, a child's weight can rebound quickly after bouts of malnutrition but height will reflect accumulated nutritional investments.

Mother's age at first birth

Mother's age at birth of first child also plays an important role in treatment heterogeneity. Figure 3.7 displays the distributional changes as mother's age at birth is increased. As evident in the graphs, the quantity-quality trade-off is amplified for mothers giving birth at a later age.

Two potential channels can explain the positive relationship between mother's age at first birth and the quantity-quality trade-off. It is possible that mothers who gave birth at an older age are better educated, and as a result, face a higher opportunity cost of providing child care themselves. Another possibility is that older mothers experience more income loss when they have an additional child compared to younger mothers. We test these two hypothesis using a simple regression:

$$Y_{ipt} = \beta_3 \widehat{nsib}_{ipt} + \alpha_1 nsib_{ipt} \times \widehat{old_mother}_i + D_{it}\gamma_4 + C_p\delta_4 + w_t + \epsilon_{ipt} \quad (3.7)$$

Equation (3.7) is the same as equation (3.3) except we include an interaction term between number of siblings and a dummy variable old_mother_i . old_mother_i takes on value 1 if mother's age at first birth is above the sample median. The dependent variables are whether the mother provides child care herself and the log of mother's income.

\widehat{nsib}_{ipt} and $nsib_{ipt} \times \widehat{old_mother}_i$ are estimated using equation (8) and equation (9), respectively.

$$nsib_{ipt} = \beta_4 coverage_i + \alpha_2 (coverage_i \times old_mother_i) + D_{it}\gamma_5 + C_p\delta_5 + w_t + v_{ipt} \quad (3.8)$$

$$\begin{aligned}
nsib_{ipt} \times old_mother_i &= \beta_5 coverage_i + \alpha_3 (coverage_i \times old_mother_i) \\
&+ D_{it} \gamma_6 + C_p \delta_6 + w_t + \eta_{ipt}
\end{aligned}
\tag{3.9}$$

Table 3.10 displays the results. Having more children increases the probability of a mother providing child care by nearly 40%. The effect is not significantly different between younger mothers and older mothers. An increase in family size decreases mother’s income by 28%. The income penalty for having more children is even larger for older mothers. Our findings suggest that a larger income penalty might be the channel behind the positive relationship between mother’s age at first birth and the quantity-quality trade-off.

The results in table 3.10 are consistent with the literature on this subject. Putz and Engelhardt (2014) find more dramatic wage losses from giving birth at a later age. They offer several explanations for this "late birth wage gap". One is that women may go through a transitional phase in the years during which later births occur.¹⁰ Another possible explanation is statistical discrimination. Putz and Engelhardt (2014) infer this since they only find a late-birth wage penalty when they use a mother’s biological age, but not when they use the number of elapsed years from joining the labor force to giving birth, as a dependent variable.

It is worth noting in the above graph that the late birth penalty is a feature of the mothers of boys as well. This effect is in stark contrast to the treatment effect heterogeneity graphs on previous pages. While there was no change in quantity-quality trade-offs for boys for varying levels of parents education there is a dramatic late-birth penalty in terms of child quality as apparent in the graph. There may need to be more research conducted in this area to examine the reasons or this.

¹⁰In analyzing Ukrainian data Nizalova et al. (2016) find that delaying birth and low education incur the largest birth-related wage penalties.

3.5 Conclusion

In this paper we use a recently developed machine learning approach to study the child quantity-quality trade-off. First we use the severity of the OCP implementation across Chinese provinces combined with the overlap of these policies with a mother's fertile years to discern exogenous variation in family size. The resulting estimation of LATE is a starting point in the examination of the heterogeneous treatment effects. We utilize a generalized random forest algorithm because it is an apt technique to detect treatment heterogeneity since it is minimizing mean square errors by contemporaneously maximizing treatment heterogeneity in the covariate space.

We detect large heterogeneous effects of family size on health-related quality measures but not on educational outcomes. Our findings suggest that parental education, income, the sex of the child and the timing of mother's first birth play the largest role in the quantity-quality trade-off. The quality penalty, especially on health outcomes, is the largest for girls. Children with more educated parents experience a smaller quality penalty. Also, an increase in mother's age at first birth induces a child quality penalty. Our findings suggest that the previous results in the literature estimating the effect of the quantity-quality trade-off may be contradictory to each other because of the instruments they were using to deal with the endogeneity of family size. As such these studies detected different complier groups and therefore different treatment effects of quantity on quality. Our approach uses a data-driven algorithm to pinpoint all sources of treatment heterogeneity. Because of this, our study can inform policy makers on precisely where to concentrate their work in their efforts to maximize child quality.

Table 3.1: Heterogeneity in the OCP implementation

Province	OCP from	Exemptions			
		Both spouses minority	One spouse minority	Both spouses aricultural	Both spouses agricultural with a girl
Beijing	1979				
Liaoning	1980	1982-1984		1985b-	1985-
Heilongjiang	1979	1981-1983 1994a-			1990-
Shanghai	1979				
Jiangsu	1979				
Shandong	1980	1984-			1986-
Henan	1981	1990c-			1990-
Hubei	1981				1988-
Hunan	1982	1990d-		1990e-	1987-
Guangxi	1982	1989a-	1985-1988a		1989-
Guizhou	1982	1982-1998		1982-1987 1988e-	1998-

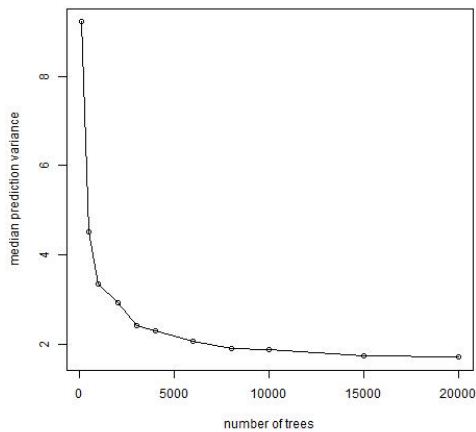
(a) only minorities with a total population less than 10 million (Manchu and Zhuang have populations exceeding 10 millions); (b) one spouse must belong to a minority, whose total population is less than 10 million; (c) both spouses must be agricultural; (d) one spouse must be agricultural; (e) one spouse must belong to a minority.

Source: Huang (2021)

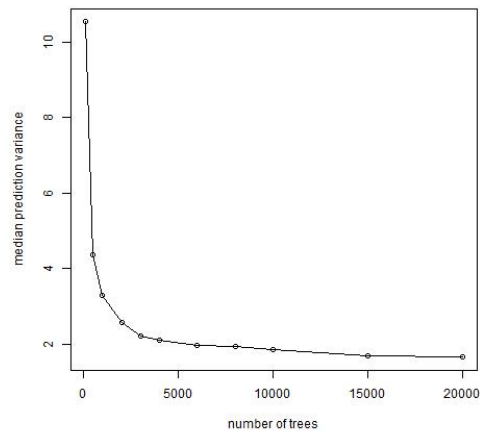
Table 3.2: Summary statistics

	All		Boys		Girls	
	Mean	Std Dev	Mean	Std Dev	Mean	Std Dev
Boy	51%	49%				
Age	11.20	3.30	11.19	3.31	11.20	3.28
HAZ	-0.42	1.35	-0.37	1.37	-0.46	1.33
WAZ	-0.41	1.33	-0.26	1.36	-0.56	1.27
Years of schooling	5.48	3.13	5.43	3.13	5.54	3.13
Currently enrolled in school	0.90	0.29	0.90	0.28	0.90	0.29
Father's height (cm)	167.73	5.79	167.71	5.76	167.74	5.82
Mother's height (cm)	156.88	5.85	157.13	5.62	156.61	6.06
Father's weight (kg)	64.91	10.06	64.93	9.84	64.89	10.29
Mother's weight (kg)	55.73	8.56	55.88	8.63	55.58	8.49
Mother's age at first birth	24.62	3.48	24.62549	3.39	24.63	3.57
Father's years of schooling	9.32	3.44	9.331241	3.43	9.30	3.45
Mother's years of schooling	8.24	4.06	8.264473	4.04	8.21	4.08

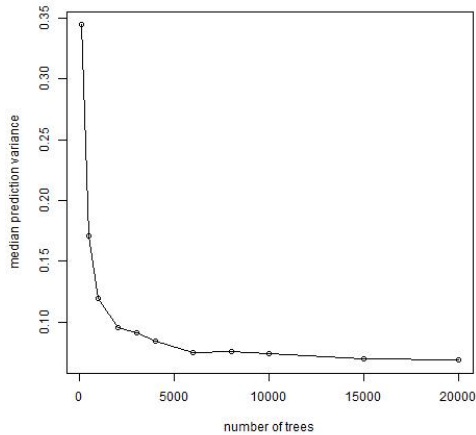
Notes: Our sample comes from the China Health and Nutrition Survey (CHNS). 8461 observations are included in the sample. HAZ is height-for-age z-score. WAZ is weight-for-age z-score.



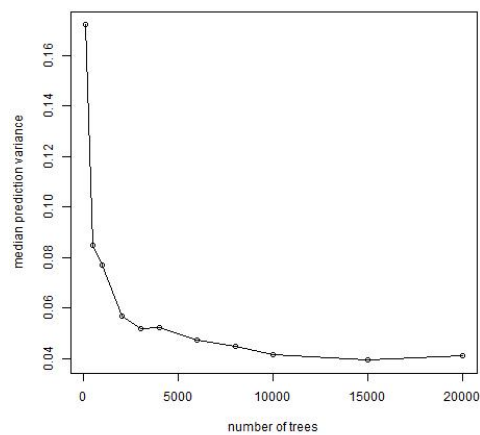
(a) HAZ median prediction variance



(b) WAZ median prediction variance



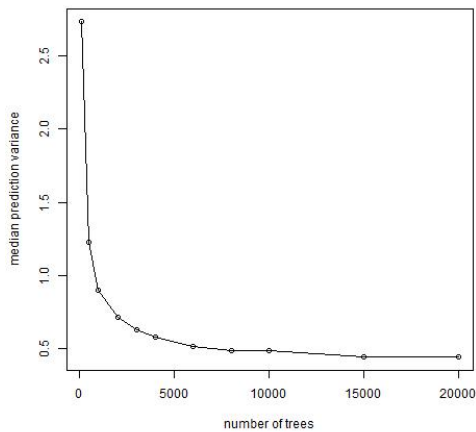
(c) relative education median prediction variance



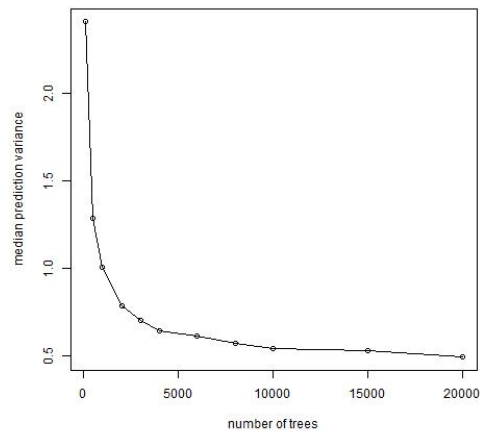
(d) school enrollment median prediction variance

Figure 3.1: Number of trees and median variance for boys

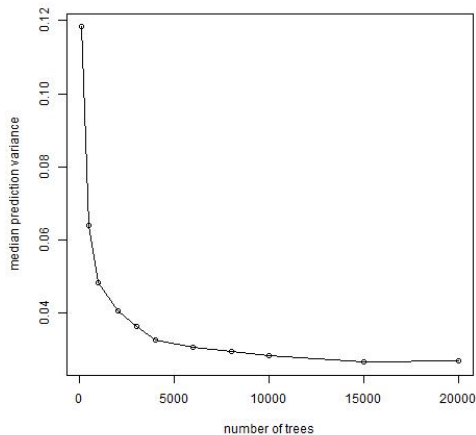
Notes: Variance of predicted treatment effect is obtained using `instrumental_forest` from R package `grf`. All the parameters are set at their default values except “num.trees”. In (a)-(d), independent variable is number of siblings; instrumental variable is *coverage*, the share of a mother’s prime fertility years covered by the strict OCP; covariate variables include a set of individual characteristic variables, a set of pre-OCP provincial characteristic variables and wave dummies. In (a) and (b), individual characteristic variables include child’s age, mom’s age at first birth, mom’s age at first birth squared, parental education-level dummies, parental weight, parental height, dummy variable for missing parental weight. In (c) and (d), individual characteristic variables include child’s age, mom’s age at first birth, mom’s age at first birth squared, parental education-level dummies. Provincial characteristic variables are the same for all panels, including sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of non-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976.



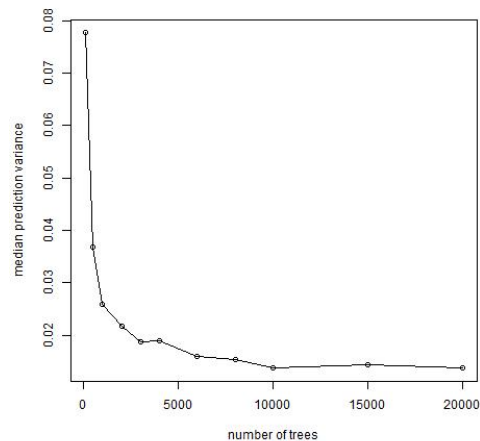
(a) HAZ median prediction variance



(b) WAZ median prediction variance



(c) relative education median prediction variance



(d) school enrollment median prediction variance

Figure 3.2: Number of trees and median variance for girls

Notes: Variance of predicted treatment effect is obtained using `instrumental_forest` from R package `grf`. All the parameters are set at their default values except “num.trees”. In (a)-(d), independent variable is number of siblings; instrumental variable is *coverage*, the share of a mother’s prime fertility years covered by the strict OCP; covariate variables include a set of individual characteristic variables, a set of pre-OCP provincial characteristic variables and wave dummies. In (a) and (b), individual characteristic variables include child’s age, mom’s age at first birth, mom’s age at first birth squared, parental education-level dummies, parental weight, parental height, dummy variable for missing parental weight, dummy variable for missing parental height. In (c) and (d), individual characteristic variables include child’s age, mom’s age at first birth, mom’s age at first birth squared, parental education-level dummies. Provincial characteristic variables are the same for all panels, including sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of non-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976.

Table 3.3: Exposure to the strict OCP and sibling sex composition

	Boys		Girls	
	(1) Male second birth	(2) Fraction of male siblings	(3) Male second birth	(4) Fraction of male siblings
<i>coverage</i>	0.0135 (0.0498)	1.799 (4.7090)	-0.123*** (0.0362)	-10.85*** (3.1830)
N	1609	1609	2095	2095

Notes: OLS regression is used in all columns. *coverage* is the share of a mother's prime fertility years covered by the strict OCP. A set of individual characteristic variables and a set of pre-OCP provincial characteristic variables are included in all the regressions. Individual characteristic variables include child's age, mom's age at first birth, mom's age at first birth squared, parental education-level dummies. Provincial characteristic variables include sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of non-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976. Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 3.4: 2SLS: Height for age z-score and weight for age z-score

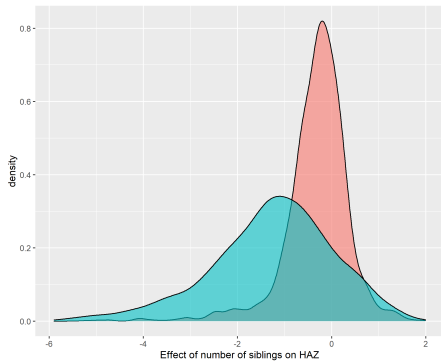
	Boys		Girls	
	(1) HAZ	(3) WAZ	(5) HAZ	(7) WAZ
Number of siblings	-1.251*** (0.363)	-1.066*** (0.337)	-0.304 (0.198)	-0.816*** (0.209)
First-stage: Number of siblings				
<i>coverage</i>	-0.186*** (0.0282)	-0.189*** (0.0282)	-0.296*** (0.0337)	-0.296*** (0.0336)
Cragg-Donald Wald F statistics	43.763	44.771	77.097	77.670
Stock-Yogo critical value: 10% maximal IV size	16.38	16.38	16.38	16.38
Individual controls	Yes	Yes	Yes	Yes
Province controls	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes
N	3648	3681	3456	3479

Notes: This table reports 2SLS regression results for HAZ and WAZ among boys and girls separately. The independent variable of interest, number of siblings, is instrumented by *coverage*, the share of a mother's prime fertility years covered by the strict OCP. Individual controls include child's age, mom's age at first birth, mom's age at first birth squared, parental education-level dummies, parental weight, parental height, dummy variable for missing parental weight, dummy variable for missing parental height. Province controls include nclude sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of non-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976. Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01.

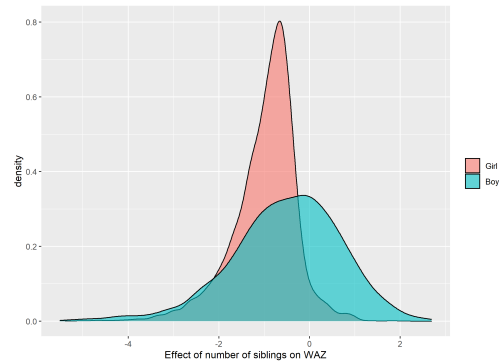
Table 3.5: 2SLS: relative education and school enrollment

	Boys		Girls	
	(5)	(7)	(6)	(8)
	Relative education	School enrollment	Relative education	School enrollment
Number of siblings	-0.158 (0.103)	-0.0480 (0.0739)	-0.149** (0.0670)	-0.0185 (0.0510)
First-stage: Number of siblings				
<i>coverage</i>	-0.197*** (0.0284)	-0.200*** (0.0277)	-0.289*** (0.0336)	-0.299*** (0.0331)
Cragg-Donald Wald F statistics	48.217	51.828	73.991	81.331
Stock-Yogo critical value: 10% maximal IV size	16.38	16.38	16.38	16.38
Individual controls	Yes	Yes	Yes	Yes
Province controls	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes
N	3648	3646	3462	3452

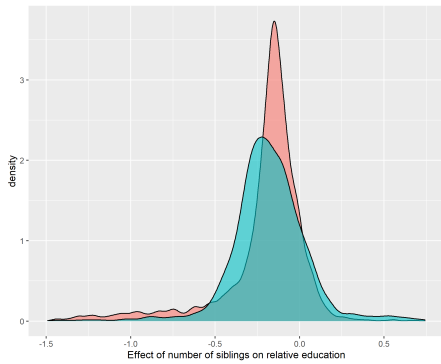
Notes: This table reports 2SLS regression results for relative education and school enrollment status among boys and girls separately. The independent variable of interest, number of siblings, is instrumented by *coverage*, the share of a mother's prime fertility years covered by the strict OCP. Individual controls include child's age, mom's age at first birth, mom's age at first birth squared, parental education-level dummies. Province controls include sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of non-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976. Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01.



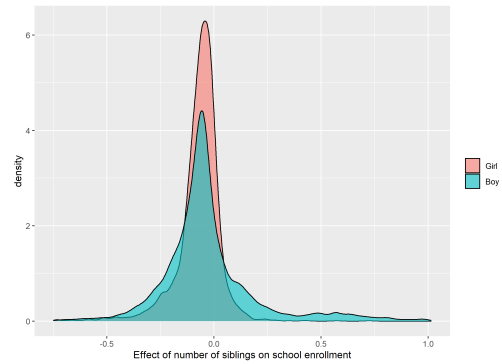
(a) Distribution of effects of number of siblings on HAZ by gender



(b) Distribution of effects of number of siblings on WAZ by gender



(c) Distribution of effects of number of siblings on relative education by gender



(d) Distribution of effects of number of siblings on relative education by gender

Figure 3.3: Distribution of effects of number of siblings by gender

Notes: Predicted treatment effect is obtained using `instrumental_forest` from R package `grf`. “num.trees” is set to be 15000. “min.node.size” is set to be 10. The rest parameters are set at their default values. In (a)-(d), independent variable is number of siblings; instrumental variable is *coverage*, the share of a mother’s prime fertility years covered by the strict OCP; covariate variables include a set of individual characteristic variables, a set of pre-OCP provincial characteristic variables and wave dummies. In (a) and (b), individual characteristic variables include child’s age, mom’s age at first birth, mom’s age at first birth squared, parental education-level dummies, parental weight, parental height, dummy variable for missing parental weight, dummy variable for missing parental height. In (c) and (d), individual characteristic variables include child’s age, mom’s age at first birth, mom’s age at first birth squared, parental education-level dummies. Provincial characteristic variables are the same for all panels, including sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of non-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976.

Table 3.6: Mean treatment effects by quartile

	HAZ		WAZ		Relative Educ		Enrolled in school	
	Boys	Girls	Boys	Girls	Boys	Girls	Boys	Girls
Mean	-1.41967	-0.26605	-0.58843	-0.98916	-0.16829	-0.21519	-0.02607	-0.07343
Mean in quartile 1 (0 - 25%)	-3.62738	-1.2622	-2.40255	-1.94055	-0.46563	-0.83078	-0.29172	-0.2085
Mean in quartile 2 (25% - 50%)	-1.55903	-0.41205	-0.83298	-1.07248	-0.23932	-0.19432	-0.0902	-0.07771
Mean in quartile 3 (50% - 75%)	-0.79682	-0.09441	-0.06766	-0.7019	-0.12462	-0.12065	-0.02931	-0.03606
Mean in quartile 4 (75% - 100%)	0.305095	0.705417	0.950052	-0.24063	0.156737	0.285109	0.30726	0.028657

Notes: Predicted treatment effect is obtained using `instrumental_forest` from R package `grf`. "num.trees" is set to be 15000. "min.node.size" is set to be 10. The rest parameters are set at their default values. In (a)-(d), independent variable is number of siblings; instrumental variable is *coverage*, the share of a mother's prime fertility years covered by the strict OCP; covariate variables include a set of individual characteristic variables, a set of pre-OCP provincial characteristic variables and wave dummies. In (a) and (b), individual characteristic variables include child's age, mom's age at first birth, mom's age at first birth squared, parental education-level dummies, parental weight, parental height, dummy variable for missing parental weight, dummy variable for missing parental height. In (c) and (d), individual characteristic variables include child's age, mom's age at first birth, mom's age at first birth squared, parental education-level dummies. Provincial characteristic variables are the same for all panels, including sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of noon-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976.

Table 3.7: Heterogeneity in local average treatment effect by predicted treatment effect

	HAZ		WAZ		Relative Educ		Enrolled in school	
	Boys	Girls	Boys	Girls	Boys	Girls	Boys	Girls
Subgroup with more negative estimated effects	-3.6584*** (0.9695)	-0.9252*** (0.3356)	-2.9751*** (0.9223)	-2.1818*** (0.6028)	-0.2392 (0.1653)	0.0600 (0.1459)	-0.0756 (0.0846)	-0.0343 (0.0867)
Rest of sample	-0.1278 (0.5285)	0.2022 (0.2649)	0.1260 (0.5680)	-0.1906 (0.2560)	-0.0706 (0.1491)	-0.2798*** (0.0782)	-0.0009 (0.1255)	-0.0262 (0.0679)
P-value, test of subgroup difference	0.0014	0.0084	0.0042	0.0024	0.4488	0.0400	0.6216	0.9415

Notes: We divide our sample into two subsamples based on the treatment effects estimated by using `instrumental_forest` from R package `grf`. Subgroup with more negative estimated effects contains individuals whose estimated treatment effects are below the median treatment effect. Rest of sample contains individuals whose estimated treatment effects are above the median treatment effect. We repeat 2SLS regressions in Table 3.4 and 3.5 using these two subgroups separately. Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 3.8: Summary statistics by predicted treatment effect on HAZ

Variable	Boys			Girls		
	Subgroup with more negative effects	Rest of sample	Difference	Subgroup with more negative effects	Rest of sample	Difference
child's age	11.2438	11.1571	0.0867	11.5695	10.9490	0.6205***
mother's age at first birth	26.0754	23.5776	2.4978***	24.8749	24.4583	0.4166***
father's height	166.5394	168.5676	-2.0282***	166.0556	168.9491	-2.8935***
mother's height	156.4396	157.6337	-1.1941***	156.6803	156.5745	0.1057
father's weight	63.4930	65.9803	-2.4873***	62.5956	66.5294	-3.9338***
mother's weight	54.6330	56.7928	-2.1598***	57.7451	54.0413	3.7038***
father's years of schooling	9.1965	9.4443	-0.2479**	8.4369	10.0430	-1.6061***
mother's years of schooling	8.2108	8.3061	-0.0953	7.0844	9.0788	-1.9944***

Notes: Subgroup with more negative effects contains individuals whose estimated treatment effects of family size on HAZ are below the mean estimated treatment effect. Rest of sample contains individuals whose estimated treatment effects of family size on HAZ are above the mean estimated treatment effect. We report covariate means for each subgroup. Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01.

Table 3.9: Summary statistics by predicted treatment effect on WAZ

Variable	WAZ			Boys			Girls		
	Subgroup with more negative effects	Rest of sample	Difference	Subgroup with more negative effects	Rest of sample	Difference	Subgroup with more negative effects	Rest of sample	Difference
child's age	11.2470	11.1542	0.0928	10.5014	11.7165	-1.2151***	10.5014	11.7165	-1.2151***
mother's age at first birth	25.7475	23.8014	1.9462***	25.2796	24.1642	1.1154***	25.2796	24.1642	1.1154***
father's height	167.1398	168.1405	-1.0006***	167.2004	168.1376	-0.9373***	167.2004	168.1376	-0.9373***
mother's height	155.8743	158.0570	-2.1826***	156.1625	156.9476	-0.7852***	156.1625	156.9476	-0.7852***
father's weight	63.9790	65.6404	-1.6613***	63.8573	65.6387	-1.7814***	63.8573	65.6387	-1.7814***
mother's weight	53.8134	57.4095	-3.5961***	55.7574	55.4570	0.3004	55.7574	55.4570	0.3004
father's years of schooling	9.0311	9.5875	-0.5564***	9.4136	9.2212	0.1924*	9.4136	9.2212	0.1924*
mother's years of schooling	8.1587	8.3478	-0.1891	8.5627	7.9445	0.6182***	8.5627	7.9445	0.6182***

Notes: Subgroup with more negative effects contains individuals whose estimated treatment effects of family size on WAZ are below the mean estimated treatment effect. Rest of sample contains individuals whose estimated treatment effects of family size on HAZ are above the mean estimated treatment effect. We report covariate means for each subgroup. Standard errors are reported in parentheses. * p < 0.1, ** p < 0.05, *** p < 0.01.

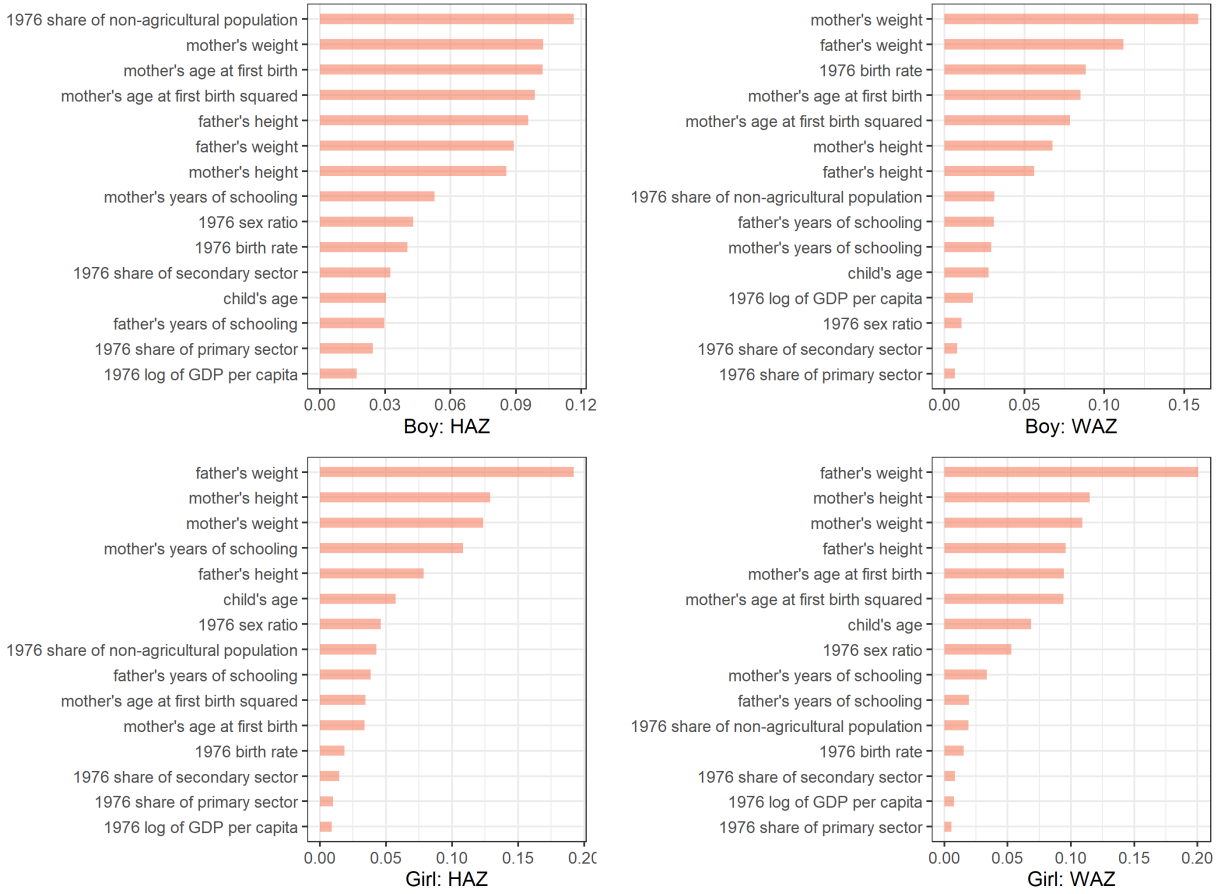
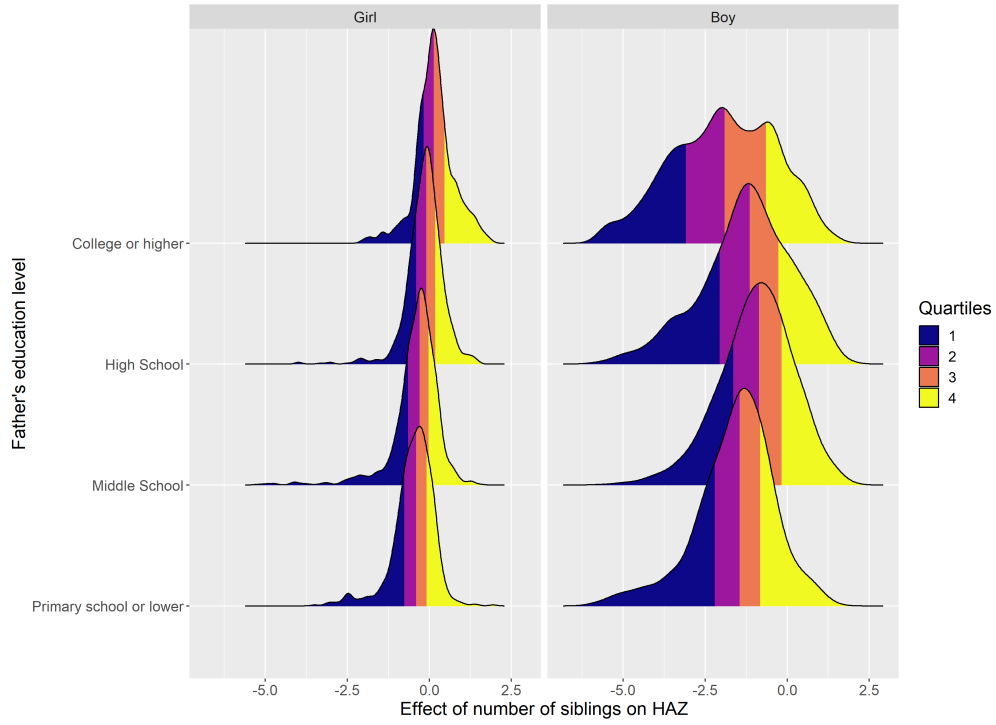
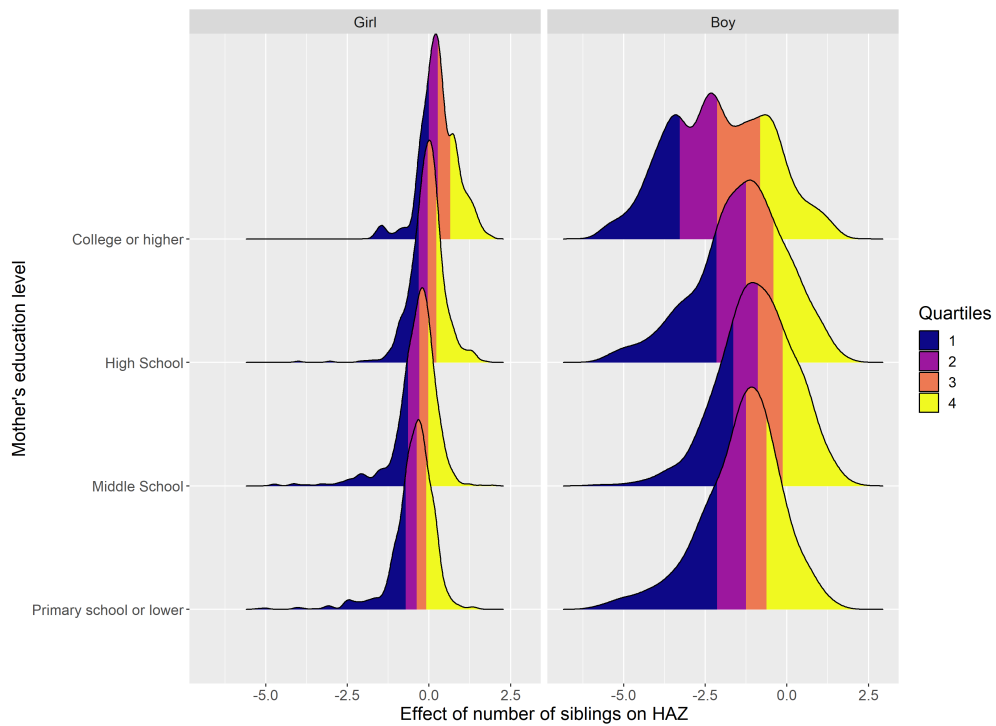


Figure 3.4: Health variable importance

Notes: Variable importance is obtained using `instrumental_forest` from R package `grf`. The variable importance measures how frequently a covariate is used in the tree splitting process.

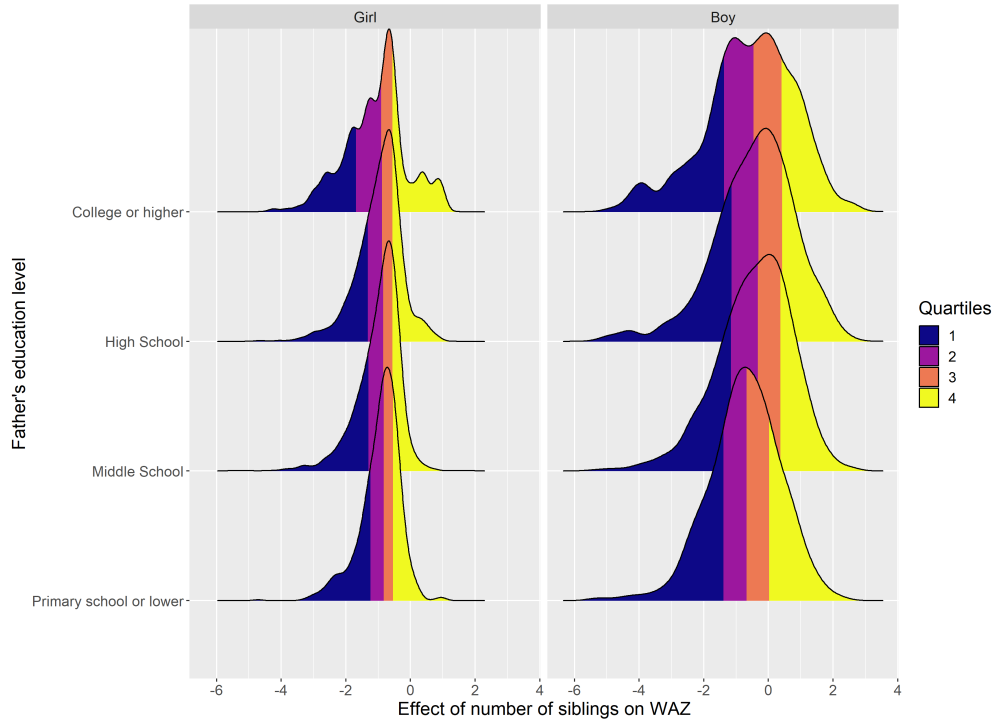


(a) by father's education

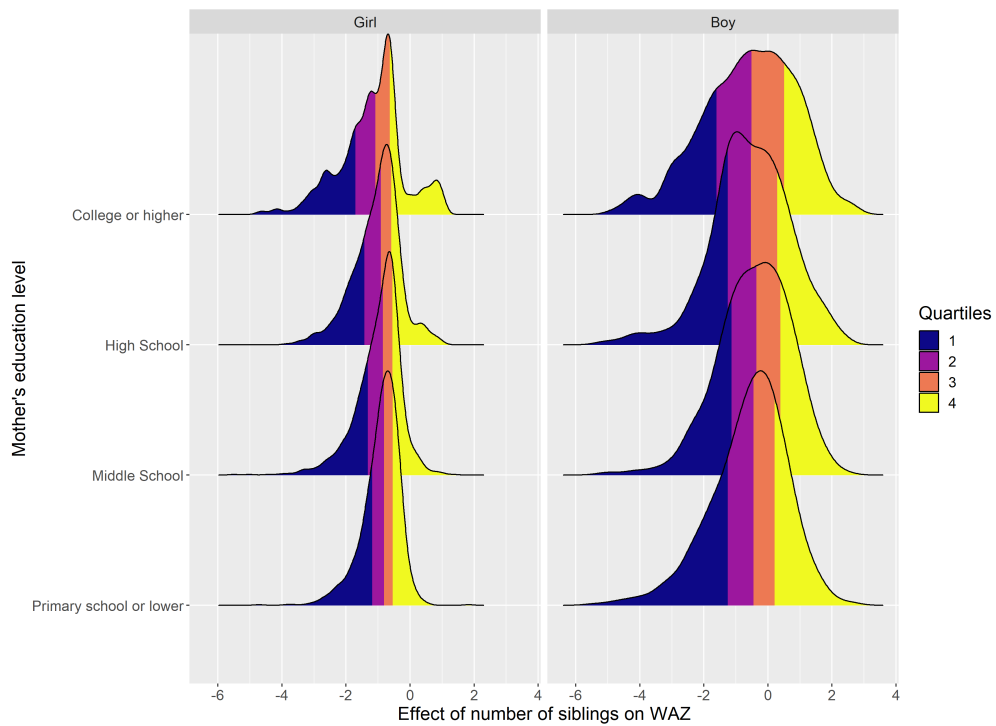


(b) by mother's education

Figure 3.5: distribution of effects of number of children on HAZ by parent's education level

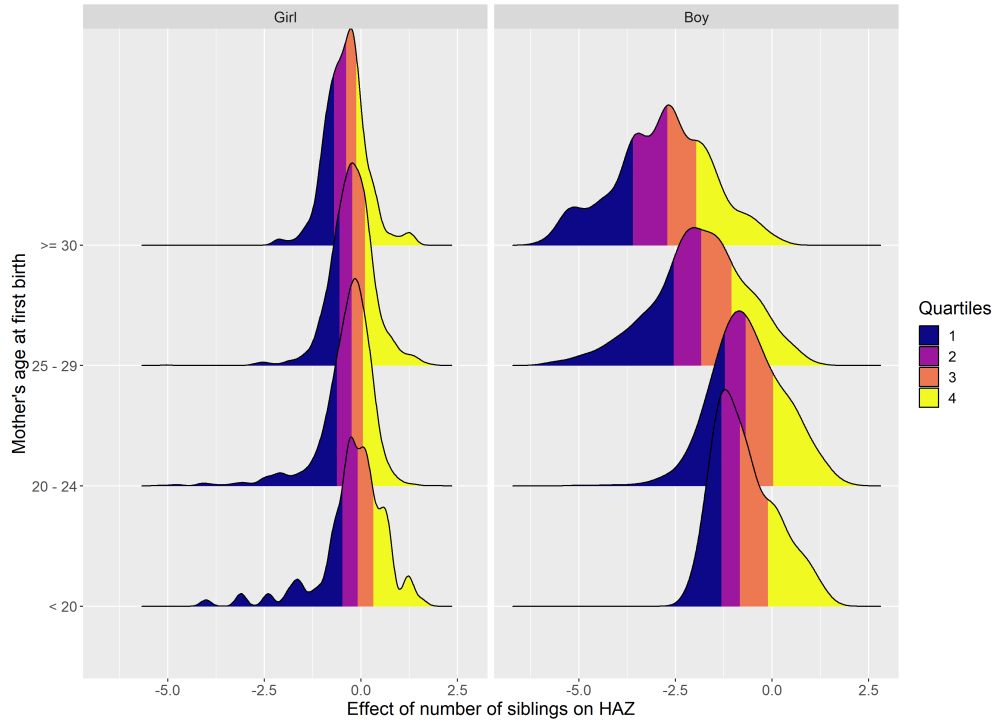


(a) by father's education

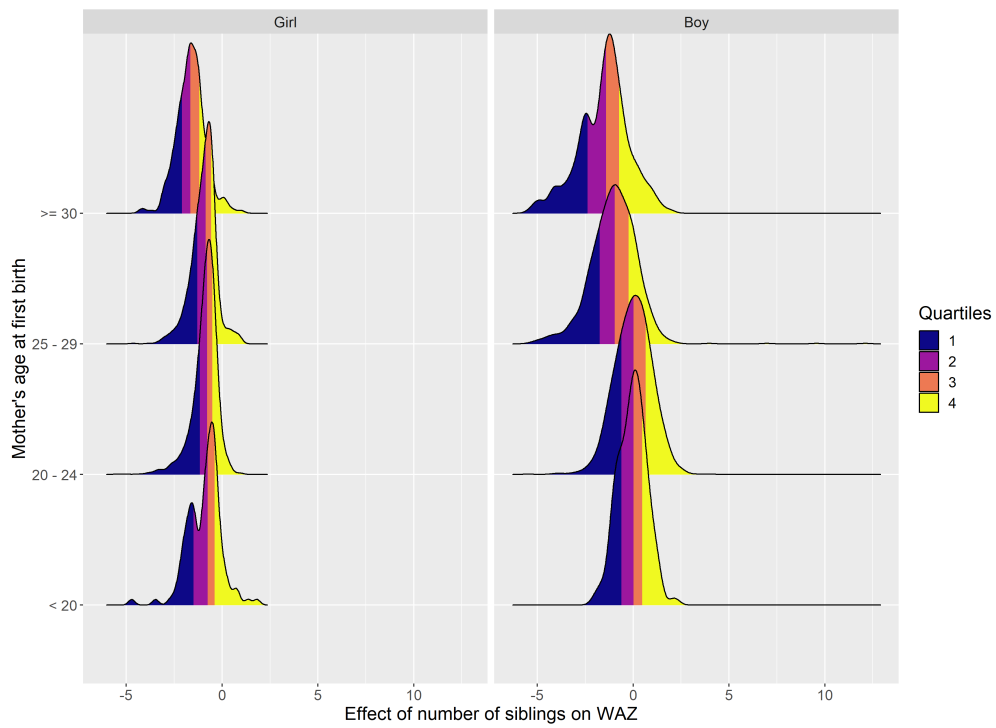


(b) by mother's education

Figure 3.6: distribution of effects of number of children on WAZ by parent's education level



(a) HAZ



(b) WAZ

Figure 3.7: distribution of effects of number of children on health by child's age group

Table 3.10: Effect of number of siblings by mother's age at first birth

	(1) child care	(2) log of mother's income
number of siblings	0.389*** (0.106)	-0.280* (0.161)
number of siblings \times old mother	-0.0270 (0.0578)	-0.292* (0.172)
N	3796	2836

Notes: Both columns report 2SLS estimates. old mother is a dummy variable that takes on value 1 if mother's age at first birth is above the sample median mother's age at first birth. Both number of siblings and number of siblings \times old mother are instrumented by *coverage* and *coverage* \times old mother. *coverage* is the share of a mother's prime fertility years covered by the strict OCP. Individual controls, province controls, and wave dummies are included in regressions of column (1) and (2). Individual controls include child's age, mom's age at first birth, mom's age at first birth squared, parental education-level dummies, parental weight, parental height, dummy variable for missing parental weight, dummy variable for missing parental height. Province controls include nclude sex ratio in 1976, birth rate in 1976, log of GDP per capita in 1976, share of non-agricultural population in 1976, share of primary industry in GDP in 1976, and share of secondary industry in GDP in 1976. Standard errors are reported in parentheses. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Bibliography

- Adkisson, R. V., and Peach, J. (1999). “Voting for president: Elections along the us-mexican border.” *Journal of Borderlands Studies*, 14(2), 67–79.
- Adkisson, R. V., and Peach, J. (2018). “The determinants of the vote for trump: an analysis of texas 2016 primary results.” *Applied Economics Letters*, 25(3), 172–175.
- Adler, P. S., and Kwon, S.-W. (2002). “Social capital: Prospects for a new concept.” *Academy of management review*, 27(1), 17–40.
- Allcott, H., and Gentzkow, M. (2017). “Social media and fake news in the 2016 election.” *Journal of economic perspectives*, 31(2), 211–36.
- Angrist, J., and Imbens, G. (1995). “Identification and estimation of local average treatment effects.”
- Angrist, J., Lavy, V., and Schlosser, A. (2010). “Multiple experiments for the causal link between the quantity and quality of children.” *Journal of Labor Economics*, 28(4), 773–824.
- Angrist, J. D., Evans, W. N., et al. (1998). “Children and their parents’ labor supply: Evidence from exogenous variation in family size.” *American Economic Review*, 88(3), 450–477.
- Athey, S., Tibshirani, J., and Wager, S. (2019). “Generalized random forests.” *The Annals of Statistics*, 47(2), 1148–1178.
- Baiardi, A., and Naghi, A. (2020). “The value added of machine learning to causal inference: Evidence from revisited studies.”

- Bailey, M., Cao, R., Kuchler, T., and Stroebel, J. (2018a). “The economic effects of social networks: Evidence from the housing market.” *Journal of Political Economy*, 126(6), 2224–2276.
- Bailey, M., Cao, R., Kuchler, T., Stroebel, J., and Wong, A. (2018b). “Social connectedness: Measurement, determinants, and effects.” *Journal of Economic Perspectives*, 32(3), 259–80.
- Ballester, C., Zenou, Y., and Calvó-Armengol, A. (2010). “Delinquent networks.” *Journal of the European Economic Association*, 8(1), 34–61.
- Bandura, A., and Walters, R. H. (1977). *Social learning theory*, vol. 1. Englewood cliffs Prentice Hall.
- Barabási, A.-L., and Albert, R. (1999). “Emergence of scaling in random networks.” *science*, 286(5439), 509–512.
- Becker, G. S. (1998). “A treatise on the family.”
- Becker, G. S., and Lewis, H. G. (1973). “On the interaction between the quantity and quality of children.” *Journal of political Economy*, 81(2, Part 2), S279–S288.
- Becker, G. S., and Tomes, N. (1976). “Child endowments and the quantity and quality of children.” *Journal of political Economy*, 84(4, Part 2), S143–S162.
- Behrman, J. R., Pollak, R. A., and Taubman, P. (1986). “Do parents favor boys?” *International Economic Review*, 33–54.
- Beugelsdijk, S., and Smulders, S. (2003). “Bridging and bonding social capital: which type is good for economic development.” In *Tilburg Anniversary Conference on Sustainable Ties in the Information Society*.
- Beugelsdijk, S., and Smulders, S. (2009). “Bonding and bridging social capital and economic growth.”
- Beyerlein, K., and Hipp, J. R. (2005). “Social capital, too much of a good thing? american religious traditions and community crime.” *Social Forces*, 84(2), 995–1013.
- Black, S. E., Devereux, P. J., and Salvanes, K. G. (2005). “The more the merrier? the

- effect of family size and birth order on children’s education.” *The Quarterly Journal of Economics*, 120(2), 669–700.
- Blandhol, C., Bonney, J., Mogstad, M., and Torgovitsky, A. (2022). “When is tsls actually late?” *University of Chicago, Becker Friedman Institute for Economics Working Paper*, (2022-16).
- Bobo, L. (1988). “Group conflict, prejudice, and the paradox of contemporary racial attitudes.” In *Eliminating racism*, 85–114, Springer.
- Bonacich, P. (1987). “Power and centrality: A family of measures.” *American journal of sociology*, 92(5), 1170–1182.
- Bonacich, P. (2007). “Some unique properties of eigenvector centrality.” *Social networks*, 29(4), 555–564.
- Bramoullé, Y., Djebbari, H., and Fortin, B. (2009). “Identification of peer effects through social networks.” *Journal of econometrics*, 150(1), 41–55.
- Breiman, L. (2001). “Random forests.” *Machine learning*, 45(1), 5–32.
- Brinch, C. N., Mogstad, M., and Wiswall, M. (2017). “Beyond late with a discrete instrument.” *Journal of Political Economy*, 125(4), 985–1039.
- Burt, R. S. (1992). *Structural holes*. Harvard university press.
- Cáceres-Delpiano, J. (2006). “The impacts of family size on investment in child quality.” *Journal of Human Resources*, 41(4), 738–754.
- Chen, Y., and Huang, Y. (2018). “The power of the government: China’s family planning leading. group and the fertility decline since 1970.” Tech. rep., GLO Discussion Paper.
- Chen, Y., and Li, H. (2009). “Mother’s education and child health: Is there a nurturing effect?” *Journal of health economics*, 28(2), 413–426.
- Chen, Y., Li, H., and Meng, L. (2013). “Prenatal sex selection and missing girls in china: Evidence from the diffusion of diagnostic ultrasound.” *Journal of Human Resources*, 48(1), 36–70.
- Coffé, H., and Geys, B. (2007). “Toward an empirical characterization of bridging and bond-

- ing social capital.” *Nonprofit and Voluntary Sector Quarterly*, 36(1), 121–139.
- Coleman, J. S. (1988). “Social capital in the creation of human capital.” *American journal of sociology*, 94, S95–S120.
- Dahl, G. B., Løken, K. V., and Mogstad, M. (2014). “Peer effects in program participation.” *American Economic Review*, 104(7), 2049–74.
- Davis, J., and Heller, S. B. (2017). “Using causal forests to predict treatment heterogeneity: An application to summer jobs.” *American Economic Review*, 107(5), 546–50.
- Eagle, N., Macy, M., and Claxton, R. (2010). “Network diversity and economic development.” *Science*, 328(5981), 1029–1031.
- Ebenstein, A. (2010). “The “missing girls” of china and the unintended consequences of the one child policy.” *Journal of Human resources*, 45(1), 87–115.
- Enke, B. (2020). “Moral values and voting.” *Journal of Political Economy*, 128(10), 3679–3729.
- Erdős, P., and Rényi, A. (1959). “Some further statistical properties of the digits in cantor’s series.” *Acta Math. Acad. Sci. Hungar.*, 10, 21–29.
- Flaxman, B. (2018). “The economic and demographic determinants of donald trump’s 2016 election victory.” *Available at SSRN 3264520*.
- Gelbach, J. B. (2002). “Public schooling for young children and maternal labor supply.” *American Economic Review*, 92(1), 307–322.
- Ghosh, R., Lerman, K., Surachawala, T., Voevodski, K., and Teng, S.-H. (2011). “Non-conservative diffusion and its application to social network analysis.” *arXiv preprint arXiv:1102.4639*.
- Glaeser, E. L., Laibson, D. I., Scheinkman, J. A., and Soutter, C. L. (2000). “Measuring trust.” *The quarterly journal of economics*, 115(3), 811–846.
- Goetz, S. J., Rupasingha, A., and Loveridge, S. (2012). “Social capital, religion, w al-m art, and hate groups in a merica.” *Social Science Quarterly*, 93(2), 379–393.
- Granovetter, M. S. (1973). “The strength of weak ties.” *American journal of sociology*, 78(6),

1360–1380.

- Guo, G., and VanWey, L. K. (1999). “Sibship size and intellectual development: Is the relationship causal?” *American Sociological Review*, 169–187.
- Hedrich, J. P. (2011). “The heterogeneous impact of the quantity-quality tradeoff: Looking beyond the mean.”
- Hotchkiss, J. L., Rupasingha, A., and Watson, T. (2022). “In-migration and dilution of community social capital.” *International Regional Science Review*, 45(1), 36–57.
- Hoyman, M., McCall, J., Paarlberg, L., and Brennan, J. (2016). “Considering the role of social capital for economic development outcomes in us counties.” *Economic Development Quarterly*, 30(4), 342–357.
- Huang, W., and Zhou, Y. (2015). “One-child policy, marriage distortion, and welfare loss.”
- Huang, Y. (2021). “Family size and children’s education: Evidence from the one-child policy in china.” *Population Research and Policy Review*, 1–26.
- Iyer, S., Kitson, M., and Toh, B. (2005). “Social capital, economic growth and regional development.” *Regional studies*, 39(8), 1015–1040.
- Jefferson, P. N., and Pryor, F. L. (1999). “On the geography of hate.” *Economics Letters*, 65(3), 389–395.
- Katz, L. (1953). “A new status index derived from sociometric analysis.” *Psychometrika*, 18(1), 39–43.
- Kessler, D. (1991). “Birth order, family size, and achievement: Family structure and wage determination.” *Journal of Labor Economics*, 9(4), 413–426.
- Kyne, D., and Aldrich, D. P. (2020). “Capturing bonding, bridging, and linking social capital through publicly available data.” *Risk, Hazards & Crisis in Public Policy*, 11(1), 61–86.
- Lee, J. (2008). “Sibling size and investment in children’s education: An asian instrument.” *Journal of Population Economics*, 21(4), 855–875.
- Lee, K., Bargagli-Stoffi, F. J., and Dominici, F. (2020). “Causal rule ensemble: Interpretable inference of heterogeneous treatment effects.” *arXiv preprint arXiv:2009.09036*.

- Li, H., Zhang, J., and Zhu, Y. (2008). “The quantity-quality trade-off of children in a developing country: Identification using chinese twins.” *Demography*, 45(1), 223–243.
- Liu, H. (2014). “The quality–quantity trade-off: evidence from the relaxation of china’s one-child policy.” *Journal of Population Economics*, 27(2), 565–602.
- Løken, K. V., Mogstad, M., and Wiswall, M. (2012). “What linear estimators miss: The effects of family income on child outcomes.” *American Economic Journal: Applied Economics*, 4(2), 1–35.
- Loury, G. C. (1976). “A dynamic theory of racial income differences.” Tech. rep., Discussion paper.
- Manski, C. F. (1993). “Identification of endogenous social effects: The reflection problem.” *The review of economic studies*, 60(3), 531–542.
- Nizalova, O. Y., Sliusarenko, T., and Shpak, S. (2016). “The motherhood wage penalty in times of transition.” *Journal of Comparative Economics*, 44(1), 56–75.
- O’Neill, E., and Weeks, M. (2018). “Causal tree estimation of heterogeneous household response to time-of-use electricity pricing schemes.” *arXiv preprint arXiv:1810.09179*.
- Ousey, G. C., and Kubrin, C. E. (2018). “Immigration and crime: Assessing a contentious issue.” *Annual Review of Criminology*, 1, 63–84.
- Overbey, L. A., Paribello, C., and Jackson, T. (2013). “Identifying influential twitter users in the 2011 egyptian revolution.” In *International Conference on Social Computing, Behavioral-Cultural Modeling, and Prediction*, 377–385, Springer.
- Putnam, R., et al. (2001). “Social capital: Measurement and consequences.” *Canadian journal of policy research*, 2(1), 41–51.
- Putnam, R. D., et al. (2000). *Bowling alone: The collapse and revival of American community*. Simon and schuster.
- Putz, T., and Engelhardt, H. (2014). “The effects of the first birth timing on women’s wages: A longitudinal analysis based on the german socio-economic panel.” *Zeitschrift für Familienforschung*, 26(3), 302–330.

- Qian, N. (2009). “Quantity-quality and the one child policy: The only-child disadvantage in school enrollment in rural china.” Tech. rep., National Bureau of Economic Research.
- Rosenzweig, M. R., and Wolpin, K. I. (1980). “Testing the quantity-quality fertility model: The use of twins as a natural experiment.” *Econometrica: journal of the Econometric Society*, 227–240.
- Rosenzweig, M. R., and Zhang, J. (2009). “Do population control policies induce more human capital investment? twins, birth weight and china’s “one-child” policy.” *The Review of Economic Studies*, 76(3), 1149–1174.
- Ruhm, C. J. (2008). “Maternal employment and adolescent development.” *Labour economics*, 15(5), 958–983.
- Rupasingha, A., Goetz, S. J., and Freshwater, D. (2006). “The production of social capital in us counties.” *The journal of socio-economics*, 35(1), 83–101.
- Sabatini, F. (2005). “Social capital as social networks. a new framework for measurement.” *A New Framework for Measurement (June 2005)*.
- Sajuria, J., vanHeerde Hudson, J., Hudson, D., Dasandi, N., and Theocharis, Y. (2015). “Tweeting alone? an analysis of bridging and bonding social capital in online networks.” *American Politics Research*, 43(4), 708–738.
- Satyanath, S., Voigtländer, N., and Voth, H.-J. (2017). “Bowling for fascism: Social capital and the rise of the nazi party.” *Journal of Political Economy*, 125(2), 478–526.
- Sunder, S., Kim, K. H., and Yorkston, E. A. (2019). “What drives herding behavior in online ratings? the role of rater experience, product portfolio, and diverging opinions.” *Journal of Marketing*, 83(6), 93–112.
- Uzzi, B. (1997). “Social structure and competition in interfirm networks...” *Administrative Science Quarterly*, 42(1), 37–69.
- Watts, D. J., and Strogatz, S. H. (1998). “Collective dynamics of ‘small-world’ networks.” *nature*, 393(6684), 440–442.
- Wilson, R. (2019). “The impact of social networks on eitc claiming behavior.” *The Review*

of Economics and Statistics, 1–45.

Wilson, R. (2021). “Isolated states of america: The impact of state borders on mobility and regional labor market adjustments.”

Woolcock, M. (1998). “Social capital and economic development: Toward a theoretical synthesis and policy framework.” *Theory and society*, 27(2), 151–208.

Zhang, J. (2017). “The evolution of china’s one-child policy and its effects on family outcomes.” *Journal of Economic Perspectives*, 31(1), 141–60.

Zhong, H. (2014). “The effect of sibling size on children’s health: a regression discontinuity design approach based on china’s one-child policy.” *China Economic Review*, 31, 156–165.