

# Codebook for Publicly Available Data on Social Capital

Raj Chetty, Matthew O. Jackson, Theresa Kuchler, Johannes Stroebel  
Nathaniel Hendren, Robert Fluegge, Sara Gong, Federico Gonzalez, Armelle Grondin,  
Matthew Jacob, Drew Johnston, Martin Koenen, Eduardo Laguna-Muggenburg,  
Florian Mudekereza, Tom Rutter, Nicolaj Thor, Wilbur Townsend, Ruby Zhang,  
Mike Bailey, Pablo Barberá, Monica Bhole, Nils Wernerfelt

July 2022

## Contents

<b>1</b>	<b>Overview</b>	<b>2</b>
<b>2</b>	<b>Codebook</b>	<b>3</b>
2.1	County-Level Data . . . . .	3
2.2	ZIP Code-Level Data . . . . .	7
2.3	High School Data . . . . .	10
2.4	College Data . . . . .	13
<b>3</b>	<b>Privacy Protection</b>	<b>16</b>
3.1	Privacy Protection for Economic Connectedness . . . . .	16
3.2	Privacy Protection for Other Statistics . . . . .	20
3.3	Minimum Cell Size Restrictions . . . . .	20
<b>4</b>	<b>Citing the Data</b>	<b>21</b>

# 1 Overview

This document presents details on the publicly available data released as part of the [Social Capital Atlas](#), constructed in [Chetty et al. \(2022a\)](#) and [Chetty et al. \(2022b\)](#). This data release includes a number of social capital statistics aggregated to the level of U.S. counties, ZIP codes, high schools, and colleges. We use methods from the differential privacy literature to add noise to these aggregate statistics to protect privacy while maintaining a high level of statistical reliability (see Section 3 for details). In addition, the data release does not include measures for cells that do not meet certain minimum size thresholds, as discussed in Section 3.3. Our measures of social capital fall into three groups:

1. **Connectedness:** The extent to which people with different characteristics (e.g., low vs. high socioeconomic status) are friends with each other.
2. **Cohesiveness:** The degree to which friendship networks are clustered into cliques and whether friendships tend to be supported by mutual friends.
3. **Civic Engagement:** Indices of participation in civic organizations or volunteering groups.

In addition to the social capital measures, in [Chetty et al. \(2022b\)](#) we decompose our measure of the degree to which people with different SES interact with each other – which we term *economic connectedness* – into two determinants: *exposure* (the extent to which people with low versus high socioeconomic status (SES) participate in the same groups) and *framing bias* conditional on exposure (the tendency for low-SES people to befriend high-SES people at lower rates even conditional on exposure). We release measures of exposure and bias as well.

As described in greater detail in [Chetty et al. \(2022a\)](#) and [Chetty et al. \(2022b\)](#), the primary analysis sample we use to construct these statistics consists of Facebook users aged between 25 and 44 who reside in the United States, were active on the Facebook platform at least once in the prior 30 days, have at least 100 U.S.-based Facebook friends, and have a non-missing residential ZIP code as of May 28, 2022. For high school and college-level statistics, we focus on individuals in the the 1986-1996 birth cohorts for measures using own SES and individuals in the 1990-2000 birth cohorts for measures using parental SES.

Note that estimates obtained from the publicly released data will not exactly match those reported in the published papers (Chetty et al. 2022a,b) because of the exclusion of small cells (e.g., those with fewer than 100 low-SES and 100 high-SES Facebook users for economic connectedness measures) and the addition of noise (see Section 3). In practice, because the estimates we report in the papers are all population-weighted, the point estimates will remain very similar but the count of the number of observations will differ. In particular, the restrictions we impose to release data publicly lead us to drop 6,034 of the 29,062 ZIP codes for which we have at least one social capital measure. These cells account for 1.25% of observations in the data, weighting by number of children with parents with below-median household income. Analogously, we exclude 3,958 of 21,483 high schools and 87 of 2,673 colleges (accounting for 2.91% and 0.03% of observations weighting by number of students, respectively), as well as 7 counties (accounting for fewer than 0.01% of observations weighting by number of children with parents with below-median household income).<sup>1</sup>

---

<sup>1</sup>For economic connectedness, the restrictions we impose lead us to drop 5,251 of 24,231 ZIP codes with non-missing economic connectedness estimates (accounting for 3.07% of observations weighting by number of children with parents with below-median household income); 9,753 of 21,361 high schools (13.41% of observations weighting by number of students); and 474 of 2,666 colleges (7.62% of observations weighting by number of students).

## 2 Codebook

### 2.1 County-Level Data

#### *County Identifiers and Population Variables*

Variable Name	Description
county	5-digit county FIPS code.
county_name	Name of the county and state.
num_below_p50	Number of children with below-national-median parental household income. This variable is not constructed using Facebook data; it is obtained from publicly available data posted at <a href="#">the Opportunity Atlas website (Chetty et al. 2018)</a> .
pop2018	Population in 2018. This variable is not constructed using Facebook data; it is obtained from publicly available data posted at <a href="#">the Census website (American Community Survey)</a> .

#### *County Connectedness Statistics*

Variable Name	Description
ec_county	Baseline definition of economic connectedness: two times the share of high-SES friends among low-SES individuals, averaged over all low-SES individuals in the county. See equations (1), (2), and (3) of <a href="#">Chetty et al. (2022a)</a> for a formal definition. We calculate SES as in Supplementary Information B.1 of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document. This variable is mapped in Figure 2A of <a href="#">Chetty et al. (2022a)</a> .
ec_se_county	The standard error of economic connectedness, incorporating both sampling error and error from the addition of noise to protect privacy. The variance due to sampling error is calculated using a bootstrap approach described in <a href="#">Chetty et al. (2022a)</a> . We then add the noise variance that we apply to protect privacy to generate a combined standard error.

<code>child_ec_county</code>	Childhood economic connectedness: two times the share of high-parental-SES friends among low-parental-SES individuals averaged over all low-parental-SES individuals in the county, calculated using only individuals' high school friends. This statistic is estimated on the subsample of individuals who can be linked to a parent with a valid SES prediction and matched to a high school. We link individuals to parents as described in Supplementary Information A.1 of Chetty et al. (2022a), and calculate SES as in Supplementary Information B.1 of Chetty et al. (2022a). When calculating childhood economic connectedness by county, we assign individuals to the counties where their high schools are located rather than counties where they currently live, in order to map people to the places where they grew up.
<code>child_ec_se_county</code>	The standard error of childhood economic connectedness, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>ec_grp_mem_county</code>	Two times the share of high-SES friends among low-SES individuals averaged over all low-SES individuals in the county, restricting attention to friendships that we can allocate to the group in which they were formed as described in Supplementary Information B.1 and B.2 of Chetty et al. (2022b).
<code>ec_high_county</code>	Economic connectedness for high-SES individuals: two times the share of high-SES friends among high-SES individuals, averaged over all high-SES individuals in the county.
<code>ec_high_se_county</code>	The standard error of economic connectedness for high-SES individuals, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>child_high_ec_county</code>	Childhood economic connectedness (calculated using only an individual's high school friends and the individual's and friends' parental SES) for high-SES individuals.
<code>child_high_ec_se_county</code>	The standard error of childhood economic connectedness for high-SES individuals, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>ec_grp_mem_high_county</code>	Two times the share of high-SES friends among high-SES individuals averaged over all high-SES individuals in the county, restricting attention to friendships that we can allocate to the group in which they were formed.
<code>exposure_grp_mem_county</code>	Mean exposure to high-SES individuals by county for low-SES individuals: two times the average share of high-SES individuals in individuals' groups, averaged over low-SES users. We assign Facebook users to groups within settings as described in Supplementary Information B.1 of Chetty et al. (2022b). We calculate SES as in Supplementary Information B.1 of Chetty et al. (2022a). This variable is plotted in Figure 4A of Chetty et al. (2022b).

<code>exposure_grp_mem_high_county</code>	Mean exposure to high-SES individuals by county for high-SES individuals: two times the average share of high-SES individuals in individuals' groups, averaged over high-SES users.
<code>child_exposure_county</code>	Mean exposure to high-parental-SES peers in high school, averaged over low-parental-SES individuals. We assign Facebook users to high schools as described in Supplementary Information B.1 of Chetty et al. (2022b). This statistic is estimated on the subsample of individuals who can be linked to a parent with a valid SES prediction and matched to a high school. We link individuals to parents as described in Supplementary Information A.1 of Chetty et al. (2022a), and calculate SES as in Supplementary Information B.1 of Chetty et al. (2022a). When calculating childhood exposure by county, we assign individuals to the counties where their high schools are located rather than counties where they currently live, in order to map people to the places where they grew up.
<code>child_high_exposure_county</code>	Mean exposure to high-parental-SES peers in high school, averaged over high-parental-SES individuals.
<code>bias_grp_mem_county</code>	<code>ec_grp_mem_county</code> divided by <code>exposure_grp_mem_county</code> , all subtracted from one. Note that this estimate of friending bias is not identical to what one would obtain by calculating friending bias at the group level and then taking means by county because of covariances between exposure and friending bias across groups (see the Exposure, bias and upward income mobility section of Methods of Chetty et al. (2022b) for further discussion). Nevertheless, these approximate estimates of friending bias are correlated above 0.85 with estimates that are aggregated up from group-level statistics. This variable is plotted in Figure 4C of Chetty et al. (2022b).
<code>bias_grp_mem_high_county</code>	<code>ec_grp_mem_high_county</code> divided by <code>exposure_grp_mem_high_county</code> , all subtracted from one.
<code>child_bias_county</code>	<code>child_ec_county</code> divided by <code>child_exposure_county</code> , all subtracted from one.
<code>child_high_bias_county</code>	<code>child_high_ec_county</code> divided by <code>child_high_exposure_county</code> , all subtracted from one.

*County Cohesiveness Statistics*

Variable Name	Description
clustering_county	The average fraction of an individual's friend pairs who are also friends with each other. See equations (4) and (5) of <a href="#">Chetty et al. (2022a)</a> . We include links to people outside the county when calculating individual clustering (equation 4), but only average clustering over individuals in the relevant county to compute clustering at the county level (equation 5). We add noise to protect privacy, as described in Section 3 of this document.
support_ratio_county	The proportion of within-county friendships where the pair of friends share a third mutual friend within the same county. See equation (6) of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document.

*County Civic Engagement Statistics*

Variable Name	Description
volunteering_rate_county	The percentage of Facebook users who are members of a group which is predicted to be about ‘volunteering’ or ‘activism’ based on group title and other group characteristics. We do not include groups that have the privacy setting ‘secret’ enabled. We additionally manually review the 50 largest such groups in the United States and the largest group in each state, and remove the very small number of groups that are clearly misclassified. We add noise to protect privacy, as described in Section 3.
civic_organizations_county	The number of Facebook Pages predicted to be “Public Good” pages based on page title, category, and other page characteristics, per 1,000 users in the county. We remove pages that do not have a website linked, do not have a description on their Facebook page or do not have an address listed. We then assign the page to a county on the basis of its listed address. We add noise to protect privacy, as described in Section 3.

## 2.2 ZIP Code-Level Data

### *ZIP Code Identifiers and Population Variables*

Variable Name	Description
zip	5-digit ZIP code tabulation area code.
county	5-digit county FIPS code.
num_below_p50	Number of children with below-national-median parental household income. This variable is not constructed using Facebook data; it is obtained from publicly available data posted at <a href="#">the Opportunity Atlas website (Chetty et al. 2018)</a> .
pop2018	Population in 2018. This variable is not constructed using Facebook data; it is obtained from publicly available data posted at <a href="#">the Census website (American Community Survey)</a> .

### *ZIP Code Connectedness Statistics*

Variable Name	Description
ec_zip	Baseline definition of economic connectedness: two times the share of high-SES friends among low-SES individuals, averaged over all low-SES individuals in the ZIP code. See equations (1), (2), and (3) of <a href="#">Chetty et al. (2022a)</a> for a formal definition. We calculate SES as in Supplementary Information B.1 of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document. This variable is mapped for the Los Angeles area in Figure 2b of <a href="#">Chetty et al. (2022a)</a> .
ec_se_zip	The standard error of economic connectedness, incorporating both sampling error and error from the addition of noise to protect privacy. The variance due to sampling error is calculated using a bootstrap approach described in <a href="#">Chetty et al. (2022a)</a> . We then add the noise variance that we apply to protect privacy to generate a combined standard error.
nbhd_ec_zip	Economic connectedness calculated using only within-neighborhood friends. We add noise to protect privacy, as described in Section 3 of this document. This variable is used to construct the green neighborhood bar in Figure 2A of <a href="#">Chetty et al. (2022b)</a> .
ec_grp_mem_zip	Two times the share of high-SES friends among low-SES individuals averaged over all low-SES individuals in the ZIP code, restricting attention to friendships that we can allocate to the group in which they were formed as described in Supplementary Information B.1 and B.2 of <a href="#">Chetty et al. (2022b)</a> . This variable is used in the first row of Table 2 of <a href="#">Chetty et al. (2022b)</a> .

<code>ec_high_zip</code>	Economic connectedness for high-SES individuals: two times the share of high-SES friends among high-SES individuals, averaged over all high-SES individuals in the ZIP code.
<code>ec_high_se_zip</code>	The standard error of economic connectedness for high-SES individuals, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>nbhd_ec_high_zip</code>	High-type economic connectedness calculated using only an individual's neighborhood friends. We add noise to protect privacy, as described in Section 3 of this document. This variable is used to construct the orange neighborhood bar in Figure 2A of Chetty et al. (2022b).
<code>ec_grp_mem_high_zip</code>	Two times the share of high-SES friends among high-SES individuals averaged over all high-SES individuals in the ZIP code, restricting attention to friendships that we can allocate to the group in which they were formed.
<code>exposure_grp_mem_zip</code>	Mean exposure to high-SES individuals by ZIP code for low-SES individuals: two times the average share of high-SES individuals in individuals' groups, averaged over low-SES users. We assign Facebook users to groups within settings as described in Supplementary Information B.1 of Chetty et al. (2022b). We calculate SES as in Supplementary Information B.1 of Chetty et al. (2022a). This variable is mapped for the Los Angeles area in Figure 4B of Chetty et al. (2022b).
<code>exposure_grp_mem_high_zip</code>	Mean exposure to high-SES individuals by ZIP code for high-SES individuals: two times the average share of high-SES individuals in individuals' groups, averaged over high-SES users.
<code>nbhd_exposure_zip</code>	Exposure calculated using only users living in the relevant ZIP code. We add noise to protect privacy, as described in Section 3 of this document. This variable is used to construct the green and orange neighborhood bars in Figure 2B of Chetty et al. (2022b). Note that this is the same for high- and low-SES individuals who live in the same ZIP code.
<code>bias_grp_mem_zip</code>	<code>ec_grp_mem_zip</code> divided by <code>exposure_grp_mem_zip</code> , all subtracted from one. This variable is mapped for the Los Angeles area in Figure 4D of Chetty et al. (2022b).
<code>bias_grp_mem_high_zip</code>	<code>ec_grp_mem_high_zip</code> divided by <code>exposure_grp_mem_high_zip</code> , all subtracted from one.
<code>nbhd_bias_zip</code>	<code>nbhd_ec_zip</code> divided by <code>nbhd_exposure_zip</code> , all subtracted from one. This variable is used to construct the green neighborhood bar in Figure 2C of Chetty et al. (2022b).

<code>nbhd_bias_high_zip</code>	<code>nbhd_ec_high_zip</code> divided by <code>nbhd_exposure_zip</code> , all subtracted from one. (Note again that, within the same neighborhood, exposure is the same for low-SES and high-SES individuals) This variable is used to construct the orange neighborhood bar in Figure 2C of <a href="#">Chetty et al. (2022b)</a> .
---------------------------------	--

### *ZIP Code Cohesiveness Statistics*

Variable Name	Description
<code>clustering_zip</code>	The average fraction of an individual's friend pairs who are also friends with each other. See equations (4) and (5) of <a href="#">Chetty et al. (2022a)</a> . We include links to people outside the ZIP code when calculating individual clustering (equation 4), but only average individual clustering over users in the relevant ZIP code to compute clustering at the ZIP code level (equation 5). We add noise to protect privacy, as described in Section 3 of this document.
<code>support_ratio_zip</code>	The proportion of within-ZIP code friendships where the pair of friends share a third mutual friend within the same ZIP code. See equation (6) of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document.

### *ZIP Code Civic Engagement Statistics*

Variable Name	Description
<code>volunteering_rate_zip</code>	The percentage of Facebook users who are members of a group which is predicted to be about 'volunteering' or 'activism' based on group title and other group characteristics. We do not include groups that have the privacy setting 'secret' enabled. We additionally manually review the 50 largest such groups in the United States and the largest group in each state, and remove the very small number of groups that are clearly misclassified. We add noise to protect privacy, as described in Section 3 of this document.
<code>civic_organizations_zip</code>	The number of Facebook Pages predicted to be "Public Good" pages based on page title, category, and other page characteristics, per 1,000 users in the ZIP code. We remove pages that do not have a website linked, do not have a description on their Facebook page or do not have an address listed. We then assign the page to a ZIP code on the basis of its listed address. We add noise to protect privacy, as described in Section 3 of this document.

## 2.3 High School Data

### *High School Identifiers and Population Variables*

Variable Name	Description
high_school	12-digit NCES school ID.
high_school_name	Name of the high school.
zip	5-digit ZIP code tabulation area code.
county	5-digit county FIPS code.
students_9_to_12	Number of students from grades 9 to 12. This variable is not constructed using Facebook data; it is obtained from publicly available data posted at <a href="#">the National Center for Education Statistics website</a>

### *High School Connectedness Statistics*

Variable Name	Description
ec_own_ses_hs	Baseline definition of economic connectedness: two times the share of high-SES friends within three birth cohorts among low-SES individuals, averaged over all low-SES individuals in the school. See equations (1), (2), and (3) of <a href="#">Chetty et al. (2022a)</a> for a formal definition. We estimate SES as in Supplementary Information B.1 of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document. This variable is used in Supplementary Information Figure 3A of <a href="#">Chetty et al. (2022b)</a> .
ec_own_ses_se_hs	The standard error of economic connectedness, incorporating both sampling error and error from the addition of noise to protect privacy. The variance due to sampling error is calculated using a bootstrap approach described in <a href="#">Chetty et al. (2022a)</a> . We then add the noise variance that we apply to protect privacy to generate a combined standard error.
ec_parent_ses_hs	Economic connectedness with parental SES: two times the share of high-parental-SES friends (who attended the same school within three birth cohorts of the individual) among low-parental-SES individuals, averaged over all low-parental-SES individuals at the school. See equations (1), (2), and (3) of <a href="#">Chetty et al. (2022a)</a> for more details on the calculation. We link individuals to parents as described in Supplementary Information A1 of <a href="#">Chetty et al. (2022a)</a> , and estimate parental SES as in Supplementary Information B.1 of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document. This variable is used in Figure 5A of <a href="#">Chetty et al. (2022b)</a> .

<code>ec_parent_ses_se_hs</code>	The standard error of economic connectedness with parental SES, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>ec_high_own_ses_hs</code>	Economic connectedness for high-SES individuals: two times the share of high-SES friends within three birth cohorts among high-SES individuals, averaged over all high-SES individuals in the school.
<code>ec_high_own_ses_se_hs</code>	The standard error of economic connectedness for high-SES individuals, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>ec_high_parent_ses_hs</code>	Economic connectedness for high-parental-SES individuals using parental SES: two times the share of high-parental-SES friends (who attended the same school within three birth cohorts of the individual) among high-parental-SES individuals, averaged over all high-parental-SES individuals in the school.
<code>ec_high_parent_ses_se_hs</code>	The standard error of economic connectedness for high-parental-SES individuals using parental SES, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>exposure_own_ses_hs</code>	Mean exposure to high-SES individuals by high school for low-SES individuals: two times the average share of high-SES individuals within three birth cohorts, averaged over low-SES users. This variable is used in Supplementary Information Figure 3A of <a href="#">Chetty et al. (2022b)</a> .
<code>exposure_parent_ses_hs</code>	Mean exposure to high-parental-SES individuals by high school for low-parental-SES individuals: two times the average share of high-parental-SES individuals within three birth cohorts, averaged over low-parental-SES users. This variable is used in Figure 5A of <a href="#">Chetty et al. (2022b)</a> .
<code>bias_own_ses_hs</code>	<code>ec_own_ses_hs</code> divided by <code>exposure_own_ses_hs</code> , all subtracted from one. This variable is used in Supplementary Information Figure 3A of <a href="#">Chetty et al. (2022b)</a> .
<code>bias_parent_ses_hs</code>	<code>ec_parent_ses_hs</code> divided by <code>exposure_parent_ses_hs</code> , all subtracted from one. This variable is used in Figure 5A of <a href="#">Chetty et al. (2022b)</a> .
<code>bias_high_own_ses_hs</code>	<code>ec_high_own_ses_hs</code> divided by <code>exposure_own_ses_hs</code> , all subtracted from one.
<code>bias_high_parent_ses_hs</code>	<code>ec_high_parent_ses_hs</code> divided by <code>exposure_parent_ses_hs</code> , all subtracted from one.

### *High School Cohesiveness Statistics*

Variable Name	Description
clustering_hs	The average fraction of an individual's friend pairs who are also friends with each other. See equations (4) and (5) of <a href="#">Chetty et al. (2022a)</a> . We include only links to friends within the school when calculating individual clustering (equation 4). We add noise to protect privacy, as described in Section 3 of this document.

### *High School Civic Engagement Statistics*

Variable Name	Description
volunteering_rate_hs	The percentage of Facebook users who are members of a group which is predicted to be about 'volunteering' or 'activism' based on group title and other group characteristics. We do not include groups that have the privacy setting 'secret' enabled. We additionally manually review the 50 largest such groups in the United States and the largest group in each state, and remove the very small number of groups that are clearly misclassified. We add noise to protect privacy, as described in Section 3.

## 2.4 College Data

### *College Identifiers and Population Variables*

Variable Name	Description
college	6-digit Office of Postsecondary Education Identification identifier (OPEID), times 100.
college_name	Name of the college.
zip	5-digit ZIP code tabulation area code.
county	5-digit county FIPS code.
mean_students_per_cohort	Mean number of students per cohort. This variable is not constructed using Facebook data; it is obtained from publicly available data posted at <a href="#">the Integrated Postsecondary Education Data System website</a> .

### *College Connectedness Statistics*

Variable Name	Description
ec_own_ses_college	Baseline definition of economic connectedness: two times the share of high-SES friends within three birth cohorts among low-SES individuals, averaged over all low-SES individuals in the college. See equations (1), (2), and (3) of <a href="#">Chetty et al. (2022a)</a> for a formal definition. We estimate SES as in Supplementary Information B.1 of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document. This variable is used in Supplementary Information Figure 3B of <a href="#">Chetty et al. (2022b)</a> .
ec_own_ses_se_college	The standard error of economic connectedness, incorporating both sampling error and error from the addition of noise to protect privacy. The variance due to sampling error is calculated using a bootstrap approach described in <a href="#">Chetty et al. (2022a)</a> . We then add the noise variance that we apply to protect privacy to generate a combined standard error.
ec_parent_ses_college	Economic connectedness with parental SES: two times the share of high-parental-SES friends (who attended the same school within three birth cohorts of the individual) among low-parental-SES individuals, averaged over all low-parental-SES individuals at the college. See equations (1), (2), and (3) of <a href="#">Chetty et al. (2022a)</a> for more details on the calculation. We link individuals to parents as described in Supplementary Information A1 of <a href="#">Chetty et al. (2022a)</a> , and estimate parental SES as in Supplementary Information B.1 of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document. This variable is used in Figure 5B of <a href="#">Chetty et al. (2022b)</a> .

<code>ec_parent_ses_se_college</code>	The standard error of economic connectedness with parental SES, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>ec_high_own_ses_college</code>	Economic connectedness for high-SES individuals: two times the share of high-SES friends within three birth cohorts among high-SES individuals, averaged over all high-SES individuals in the college.
<code>ec_high_own_ses_se_college</code>	The standard error of economic connectedness for high-SES individuals, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>ec_high_parent_ses_college</code>	Economic connectedness for high-parental-SES individuals using parental SES: two times the share of high-parental-SES friends (who attended the same school within three birth cohorts of the individual) among high-parental-SES individuals, averaged over all high-parental-SES individuals in the college.
<code>ec_high_parent_ses_se_college</code>	The standard error of economic connectedness for high-parental-SES individuals using parental SES, incorporating both sampling error and error from the addition of noise to protect privacy.
<code>exposure_own_ses_college</code>	Mean exposure to high-SES individuals by college for low-SES individuals: two times the average share of high-SES individuals within three birth cohorts, averaged over low-SES users. This variable is used in Supplementary Information Figure 3B of <a href="#">Chetty et al. (2022b)</a> .
<code>exposure_parent_ses_college</code>	Mean exposure to high-parental-SES individuals by college for low-parental-SES individuals: two times the average share of high-parental-SES individuals within three birth cohorts, averaged over low-parental-SES users. This variable is used in Figure 5B of <a href="#">Chetty et al. (2022b)</a> .
<code>bias_own_ses_college</code>	<code>ec_own_ses_college</code> divided by <code>exposure_own_ses_college</code> , all subtracted from one. This variable is used in Supplementary Information Figure 3B of <a href="#">Chetty et al. (2022b)</a> .
<code>bias_parent_ses_college</code>	<code>ec_parent_ses_college</code> divided by <code>exposure_parent_ses_college</code> , all subtracted from one.
<code>bias_high_own_ses_college</code>	<code>ec_high_own_ses_college</code> divided by <code>exposure_own_ses_college</code> , all subtracted from one.
<code>bias_high_parent_ses_college</code>	<code>ec_high_parent_ses_college</code> divided by <code>exposure_parent_ses_college</code> , all subtracted from one.

*College Cohesiveness Statistics*

Variable Name	Description
<code>clustering_college</code>	The average fraction of an individual's friend pairs who are also friends with each other. See equations (4) and (5) of <a href="#">Chetty et al. (2022a)</a> . We include only links to friends within the college when calculating individual clustering (equation 4). We add noise to protect privacy, as described in Section 3 of this document.
<code>support_ratio_college</code>	The proportion of within-college friendships where the pair of friends share a third mutual friend within the same college. See equation (6) of <a href="#">Chetty et al. (2022a)</a> . We add noise to protect privacy, as described in Section 3 of this document.

*College Civic Engagement Statistics*

Variable Name	Description
<code>volunteering_rate_college</code>	The percentage of Facebook users who are members of a group which is predicted to be about 'volunteering' or 'activism' based on group title and other group characteristics. We do not include groups that have the privacy setting 'secret' enabled. We additionally manually review the 50 largest such groups in the United States and the largest group in each state, and remove the very small number of groups that are clearly misclassified. We add noise to protect privacy, as described in Section 3.

### 3 Privacy Protection

Privacy in these datasets are protected using tools from differential privacy, which adds enough noise to the data to provide precise guarantees that no significant additional information can be learned from the data about individuals (beyond what is already available from any external source). The objective of our approach to privacy protection is to add sufficient noise and to aggregate over sufficiently many individuals that it is not reasonably possible to learn about any one individual from the data. In particular, when releasing an aggregated statistic calculated over hundreds of individuals, it is possible to provide strong privacy protections in the above sense while applying only a small amount of noise, since each individual's measurement contributes a small amount to the overall statistic. Thus, the statistical value of the data is maintained while the privacy of any one individual is protected.

This section describes the methods we use to protect privacy, adapting techniques developed in [Chetty and Friedman \(2019\)](#) to release the statistics on upward mobility in [Chetty et al. \(2018\)](#) to our network setting. We thank Salil Vadhan for his help in developing the methods used below.

#### 3.1 Privacy Protection for Economic Connectedness

##### 3.1.1 Notation

- $i$  and  $j$  denote a generic pair of individuals in the primary analysis sample.
- Let  $c$  be a cell: the set of Facebook users whom we are considering, such as students in a certain school or residents of a county; let  $N_c$  be the number of users in cell  $c$ .
- Let  $\mathbf{g} \in \{0, 1\}^{n \times n}$  be a matrix representing the friendships we are considering for the users in cell  $c$ . In some cases (such as schools), we only consider friendships within the same school, so  $n = N_c$ . In other cases, such as county, we also consider friendships outside the county, so  $n > N_c$ .
- Let  $L$  denote the low-SES agents and  $H$  denote the high-SES agents.
- Let  $N_{Lc}$  denote the number of low-SES users in cell  $c$ .
- Let  $\mathbf{g} - j$  denote the subgraph induced on the network  $\mathbf{g}$  by removing node  $j$  and all related edges from the network.
- Let  $\mathbf{g} + j$  denote a new network in which a node  $j$  has been added, together with some edges.
- Let  $d_i(\mathbf{g}) = \sum_j g_{ij}$  denote the number of friends  $i$  has ( $i$ 's degree), and  $H_i(\mathbf{g}) = \sum_{j \in H} g_{ij}$  denote the number of high-SES friends  $i$  has.

##### 3.1.2 Economic Connectedness

We calculate economic connectedness ( $EC_c$ ) of  $c$  as follows:

$$\begin{aligned} IEC_i(\mathbf{g}) &\equiv \left\{ \frac{H_i(\mathbf{g})}{d_i(\mathbf{g})} \right\} / 0.5 \\ EC_c(\mathbf{g}) &= \frac{\sum_{i \in L \cap c} IEC_i(\mathbf{g})}{N_{Lc}} \end{aligned}$$

### 3.1.3 Calculating Local Sensitivity

We characterize the local sensitivity of EC with respect to both deletions and additions of a node. Cells always have more than one user and their degrees are always more than one, so we never have to worry about division by 0 in what follows. The local sensitivity of  $EC_c(\mathbf{g})$  is defined as:

$$LS_c(\mathbf{g}) \equiv \max \left[ \max_j |EC_c(\mathbf{g} + j) - EC_c(\mathbf{g})|, \max_j |EC_c(\mathbf{g}) - EC_c(\mathbf{g} - j)| \right].$$

In what follows, the terms  $H_i$  and  $d_i$  are always relative to the starting network  $\mathbf{g}$  and so we omit that notation. The following result characterizes the local sensitivity of Economic Connectedness as a function of the network, and determines how much noise we apply to the raw connectedness statistic in each cell.

**Theorem 1.** *The local sensitivity of EC for cell  $c$  at a given network  $\mathbf{g}$  is:*

$$LS_c(\mathbf{g}) = \max \left\{ \frac{2}{N_{Lc}} \sum_{i \in L \cap c} \frac{d_i - H_i}{d_i(d_i - 1)}, \frac{2}{N_{Lc} - 1} \sum_{i \in L \cap c} \frac{H_i}{d_i(d_i - 1)}, \frac{2}{N_{Lc}} \right\}$$

*Proof.* The proof proceeds by considering four scenarios.

1. The addition or deletion of a high-SES node that has negative influence.
2. The addition or deletion of a high-SES node that has positive influence.
3. The addition or deletion of a low-SES node that has negative influence.
4. The addition or deletion of a low-SES node that has positive influence.

In each case, we consider the effect of an arbitrary addition and then an arbitrary deletion, showing that an arbitrary deletion has a larger bound, and hence is the relevant term for computation of sensitivity.

**Case 1:** A high-SES node can only have (weakly) positive influence because it does not directly enter the EC sum and can only cause the IEC of terms in the sum to increase; hence, this case is not relevant.

**Case 2 Additions:** The most an *additional* high-SES node can move EC up by occurs when a high-SES node enters which befriends every low-SES node. Then the change in EC is:

$$\begin{aligned} & \frac{1}{N_{Lc}} \left( 2 \sum_{i \in L \cap c} \frac{H_i + 1}{d_i + 1} - 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} \right) \\ &= \frac{1}{N_{Lc}} \left( 2 \sum_{i \in L \cap c} \frac{d_i - H_i}{d_i(d_i + 1)} \right) \end{aligned}$$

**Case 2 Deletions:** The most the *removal* of high-SES node can move EC down occurs when one removes a high-SES node who was friends with every low-SES node. Then the change in EC is:

$$\begin{aligned} & \frac{1}{N_{Lc}} \left( 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} - 2 \sum_{i \in L \cap c} \frac{H_i - 1}{d_i - 1} \right) \\ &= \frac{1}{N_{Lc}} \left( 2 \sum_{i \in L \cap c} \frac{d_i - H_i}{d_i(d_i - 1)} \right) \end{aligned}$$

Note that this is larger than the case of additions (the denominator for each term in the sum is smaller for deletions), and is the first term of the max operator of Theorem 1.

**Case 3 Additions:** The most an *additional* low-SES node can move EC down is when the IEC of the new node is 0 and it befriends every other low-SES node.

$$\begin{aligned} & \frac{1}{N_{Lc}} \left( 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} \right) - \frac{1}{N_{Lc} + 1} \left( 0 + 2 \sum_{i \in L \cap c} \frac{H_i}{d_i + 1} \right) \\ & \leq \frac{1}{N_{Lc} + 1} \left( 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} - 0 - 2 \sum_{i \in L \cap c} \frac{H_i}{d_i + 1} \right) \\ & = \frac{1}{N_{Lc} + 1} \left( 2 \sum_{i \in L \cap c} \frac{H_i}{d_i(d_i + 1)} \right) \end{aligned}$$

**Case 3 Deletions:** The most the *removal* of a low-SES node can move EC up is when the IEC of the deleted node is 0 and it was friends with every other low-SES node.

$$\begin{aligned} & \frac{1}{N_{Lc} - 1} \left( \sum_{i \in L \cap c, i \neq j} \frac{H_i}{d_i - 1} - 0 \right) - \frac{1}{N_{Lc}} \sum_{i \in L \cap c} \frac{H_i}{d_i} \\ & \leq \frac{1}{N_{Lc} - 1} \left( \sum_{i \in L \cap c, i \neq j} \frac{H_i}{d_i - 1} - 0 - 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} \right) \\ & \leq \frac{1}{N_{Lc} - 1} \left( 2 \sum_{i \in L \cap c} \frac{H_i}{d_i - 1} - 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} \right) \\ & = \frac{1}{N_{Lc} - 1} \left( 2 \sum_{i \in L \cap c} \frac{H_i}{d_i(d_i - 1)} \right) \end{aligned}$$

Note that this is larger than the case of additions, and is the second term of the max operator of Theorem 1.

**Case 4 Additions:** The most an *additional* low-SES node can move EC up is when the additional low-SES node only befriends high-SES nodes. Then it does not change the IEC of any other low-SES node down and its own IEC is maximized.

$$\begin{aligned} & \frac{1}{N_{Lc} + 1} \left( 2 + 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} \right) - \frac{2}{N_{Lc}} \sum_{i \in L \cap c} \frac{H_i}{d_i} \\ & \leq \frac{1}{N_{Lc}} \left( 2 + 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} - 2 \sum_{i \in L \cap c} \frac{H_i}{d_i} \right) \\ & = \frac{2}{N_{Lc}} \end{aligned}$$

**Case 4 Deletions:** The most a *removal* of a low-SES node can move EC down is when the removed low-SES node was only friends with high-SES nodes. Then its own IEC was maximized

and it did not bring down the IEC of other low-SES nodes.

$$\begin{aligned}
& \frac{2}{N_{Lc}} \sum_{i \in L \cap c} \frac{H_i}{d_i} - \frac{2}{N_{Lc} - 1} \left( \sum_{i \in L \cap c, i \neq j} \frac{H_i}{d_i} \right) \\
& \leq \frac{2}{N_{Lc}} \left( 1 + \sum_{i \in L \cap c, i \neq j} \frac{H_i}{d_i} \right) - \frac{2}{N_{Lc} - 1} \left( \sum_{i \in L \cap c, i \neq j} \frac{H_i}{d_i} \right) \\
& \leq \frac{2}{N_{Lc}} + \frac{2}{N_{Lc} - 1} \left( \sum_{i \in L \cap c, i \neq j} \frac{H_i}{d_i} \right) - \frac{2}{N_{Lc} - 1} \left( \sum_{i \in L \cap c, i \neq j} \frac{H_i}{d_i} \right) \\
& = \frac{2}{N_{Lc}}
\end{aligned}$$

This is the third term of the max operator.

Theorem 1 follows from combining the deletion cases for cases 2, 3, and 4.  $\square$

### 3.1.4 Constructing an Envelope from the Sensitivities

Using the local sensitivities  $S_c$  for each cell, we construct a smooth envelope based on one non-noisy parameter  $\chi$  which we do not release to the public.

$$\chi = \max_c \left\{ \frac{S_c}{\frac{1}{N_{Lc}} \sum_{i \in L \cap c} \frac{1}{d_i}} \right\}$$

We then calculate smoothed noise  $\tilde{S}_c$  as:

$$\tilde{S}_c = \chi \times \frac{1}{N_{Lc}} \sum_{i \in L \cap c} \frac{1}{d_i}$$

applying noise to EC calculated in each cell  $c$  from the distribution:

$$\text{Laplace} \left( 0, \frac{\tilde{S}_c}{\varepsilon} \right)$$

using  $\varepsilon = 8$ , as in Chetty et al. (2018).

### 3.2 Privacy Protection for Other Statistics

We protect privacy for the other statistics we release as follows.

**Exposure and Volunteering Rate.** These variables are simple means over independent values. Individual-level values of exposure lie between 0 and 2, so we follow standard results in the differential privacy literature for means of bounded variables and apply noise from the  $\text{Laplace}(0, 2/N\epsilon)$  distribution, where  $N$  is the number of users in the cell. Individual-level volunteering is a binary value equal to either zero or one, so we apply noise from the  $\text{Laplace}(0, 1/N\epsilon)$  distribution.

**Friending Bias.** We approximate friending bias as the ratio of two privacy-protected statistics (EC and exposure) we release publicly; since it is simply a function of publicly available, privacy-protected statistics, no further noise is added.

**Clustering and Support Ratio.** For cohesiveness measures (clustering and support ratio), which do not use any information on individuals' characteristics, we follow the privacy procedures developed for the Social Connectedness Index, which was [released](#) in Fall 2020 by the Facebook Data for Good team (see [Bailey et al., 2018, 2020, 2021](#)). Specifically, we compute the statistic over the subgraph from a 99% random sample of users. We then apply additional noise from the  $\text{Laplace}(0, 0.001/8)$  distribution to the cell-level averages of the node-level network statistics.

**Civic Organizations.** Public good page density is a variable based on a count (the number of pages in an area) and only indirectly on users (through the density calculation), so we add noise from the  $\text{Laplace}(0, 0.001/8)$  distribution.

### 3.3 Minimum Cell Size Restrictions

Finally, to further protect privacy, we only release statistics on economic connectedness, exposure, and friending bias for cells that contain at least 100 low-SES and at least 100 high-SES Facebook users. We only release statistics on volunteering rates, clustering, and support ratios for cells that contain at least 100 Facebook users.

## 4 Citing the Data

Please cite the following two publications as the source of the data:

- Chetty, Raj, Matthew O. Jackson, Theresa Kuchler, Johannes Stroebel, Nathaniel Hendren, Robert Fluegge, Sara Gong, Federico Gonzalez, Armelle Grondin, Matthew Jacob, Drew Johnston, Martin Koenen, Eduardo Laguna-Muggenberg, Florian Mudekereza, Tom Rutter, Nicolaj Thor, Wilbur Townsend, Ruby Zhang, Mike Bailey, Pablo Barberá, Monica Bhole, and Nils Wernerfelt (2022a). “Social Capital I: Measurement and Associations with Economic Mobility.” *Nature*, 608(7921), 108–121.

```
@article{  
  chetty2022socialcapitalone,  
  title = {Social Capital I: Measurement and Associations with Economic Mobility},  
  author = {  
    Chetty, Raj and Jackson, Matthew O. and Kuchler, Theresa and  
    Stroebel, Johannes and Hendren, Nathaniel and Fluegge, Robert and  
    Gong, Sara and Gonzalez, Federico and Grondin, Armelle and  
    Jacob, Matthew and Johnston, Drew and Koenen, Martin and  
    Laguna-Muggenberg, Eduardo and Mudekereza, Florian and Rutter, Tom and  
    Thor, Nicolaj and Townsend, Wilbur and Zhang, Ruby and  
    Bailey, Mike and Barber\'{a}, Pablo and Bhole, Monica and Wernerfelt, Nils},  
  journal = {Nature},  
  volume = {608},  
  number = {7921},  
  pages = {108$-$121},  
  year = {2022}  
}
```

- Chetty, Raj, Matthew O. Jackson, Theresa Kuchler, Johannes Stroebel, Nathaniel Hendren, Robert Fluegge, Sara Gong, Federico Gonzalez, Armelle Grondin, Matthew Jacob, Drew Johnston, Martin Koenen, Eduardo Laguna-Muggenberg, Florian Mudekereza, Tom Rutter, Nicolaj Thor, Wilbur Townsend, Ruby Zhang, Mike Bailey, Pablo Barberá, Monica Bhole, and Nils Wernerfelt (2022b). “Social Capital II: Determinants of Economic Connectedness.” *Nature*, 608(7921), 122–134.

```
@article{  
  chetty2022socialcapitaltwo,  
  title = {Social Capital II: Determinants of Economic Connectedness},  
  author = {  
    Chetty, Raj and Jackson, Matthew O. and Kuchler, Theresa and  
    Stroebel, Johannes and Hendren, Nathaniel and Fluegge, Robert and  
    Gong, Sara and Gonzalez, Federico and Grondin, Armelle and  
    Jacob, Matthew and Johnston, Drew and Koenen, Martin and  
    Laguna-Muggenberg, Eduardo and Mudekereza, Florian and Rutter, Tom and  
    Thor, Nicolaj and Townsend, Wilbur and Zhang, Ruby and  
    Bailey, Mike and Barber\'{a}, Pablo and Bhole, Monica and Wernerfelt, Nils},  
  journal = {Nature},  
  volume = {608},  
  number = {7921},  
  pages = {122$-$134},  
  year = {2022}  
}
```

## References

- Bailey, M., R. Cao, T. Kuchler, J. Stroebel, and A. Wong (2018). Social connectedness: Measurement, determinants, and effects. *Journal of Economic Perspectives* 32(3), 259–80.
- Bailey, M., A. Gupta, S. Hillenbrand, T. Kuchler, R. Richmond, and J. Stroebel (2021). International trade and social connectedness. *Journal of International Economics* 129, 103418.
- Bailey, M., D. Johnston, T. Kuchler, D. Russel, J. Stroebel, et al. (2020). The determinants of social connectedness in Europe. In *International Conference on Social Informatics*, pp. 1–14. Springer.
- Chetty, R. and J. N. Friedman (2019, October). A practical method to reduce privacy loss when disclosing statistics based on small samples. *Journal of Privacy and Confidentiality* 9(2).
- Chetty, R., J. N. Friedman, N. Hendren, M. R. Jones, and S. R. Porter (2018). The Opportunity Atlas: Mapping the Childhood Roots of Social Mobility. Working Paper 25147, National Bureau of Economic Research.
- Chetty, R., M. O. Jackson, T. Kuchler, J. Stroebel, N. Hendren, R. Fluegge, S. Gong, F. Gonzalez, A. Grondin, M. Jacob, D. Johnston, M. Koenen, E. Laguna-Muggenberg, F. Mudekereza, T. Rutter, N. Thor, W. Townsend, R. Zhang, M. Bailey, P. Barberá, M. Bhole, and N. Wernerfelt (2022a). Social capital I: Measurement and associations with economic mobility. *Nature* 608(7921), 108–121.
- Chetty, R., M. O. Jackson, T. Kuchler, J. Stroebel, N. Hendren, R. Fluegge, S. Gong, F. Gonzalez, A. Grondin, M. Jacob, D. Johnston, M. Koenen, E. Laguna-Muggenberg, F. Mudekereza, T. Rutter, N. Thor, W. Townsend, R. Zhang, M. Bailey, P. Barberá, M. Bhole, and N. Wernerfelt (2022b). Social capital II: Determinants of economic connectedness. *Nature* 608(7921), 122–134.