

## Research Article

Amy Tang, Yang Wang, Chunqiang Tang\*

# Understanding the Factors Influencing the Number of Extracurricular Clubs in American High Schools

<https://doi.org/10.1515/edu-2025-0093>

received December 21, 2024; accepted May 28, 2025

**Abstract:** Previous research has provided compelling evidence of a strong connection between extracurricular activities and positive youth development. While both school offerings and student participation affect the outcomes of extracurricular activities, earlier studies have primarily focused on student participation. In contrast, this study shifts the focus to school offerings. We extensively collect lists of extracurricular clubs offered by hundreds of American high schools and analyze the relationship between schools' club counts and various factors, such as school enrollment, household income, pupil-to-teacher ratio, and racial demographics. We find that, although the relationship exhibits complex higher-order effects and nonlinearity, it can still be effectively captured by our carefully constructed predictive model. Moreover, we find that, despite the significant influence of school demographics, schools still have ample opportunities to take the initiative to improve their club offerings.

## 1 Introduction

American Youths have been active in extracurricular activities (Mahoney, Harris, & Eccles, 2006), and the past studies have provided compelling evidence of a strong connection between positive youth development and active engagement in organized extracurricular activities (Buckley & Lee, 2021; Feraco, Resnati, Fregonese, Spoto, & Meneghetti, 2023; Heath, Anderson, Turner, & Payne, 2022; Mkude & Mubofu, 2022; Leksuwanakun, Dangprapai, & Wangsaturaka,

2023; Busseri, Rose-Krasnor, Willoughby, & Chalmers, 2006; Gilman, Meyers, & Perez, 2004; Eccles, Barber, Stone, & Hunt, 2003; Marsh & Kleitman, 2002; Feldman & Matjasko, 2005; Darling, Caldwell, & Smith, 2005; Eccles & Templeton, 2002; Gardner, Roth, J., & Brooks-Gunn, 2008; Fredricks & Eccles, 2006; Reeves, 2008; Fredricks & Eccles, 2006; Fredricks & Eccles, 2005; Peck, Roeser, Zarrett, & Eccles, 2008). Participation in such activities has been found to be linked to a reduction in problem behaviors and an improvement in academic performance (Mahoney, Larson, & Eccles, 2005; Anjum, 2021; Mukesh, Acharya, & Pillai, 2023). Furthermore, these advantages may extend into young adulthood, predicting academic achievements and fostering prosocial behaviors (Zaff, Moore, Papillo, & Williams, 2003). More recently, extracurricular activities have also been found to help prevent social networks addiction (Borrego & Cuadrado, 2025).

The benefits of high school extracurricular clubs are rooted in several foundational learning theories. Interest-based learning theory (Dewey & Wheeler, 1913) asserts that students learn most effectively when education aligns with their interests and incorporates active, hands-on experiences. Self-determination theory (Deci & Ryan, 2012) suggests that individuals are motivated to learn and grow when they experience autonomy, competence, and relatedness. Social learning theory (Bandura & Walters, 1977) underscores the role of observation and imitation in the learning process, while the situated learning theory (Lave, 1991) emphasizes that learning occurs through participation in communities of practice, engaging in authentic activities. Finally, the expectancy-value theory (Wigfield & Eccles, 2000) highlights that individuals' beliefs about their abilities and the value they assign to a task influence their motivation and performance. Extracurricular clubs align closely with these theories by catering to students' interests, fostering autonomy in self-organization, providing opportunities to observe, practice, and collaborate within a community, and enhancing confidence through successful experiences in activities that match their interests.

Given the importance of extracurricular clubs in student development, researchers have studied factors influencing

\* **Corresponding author: Chunqiang Tang**, Meta platforms, San Jose, United States, e-mail: tangchq@gmail.com

**Amy Tang:** Lynbrook High School, San Jose, United States

**Yang Wang:** Department of Computer Science and Engineering, The Ohio State University, Ohio, United States

participation for over half a century. However, significant gaps remain, which this study seeks to address.

First, previous studies have primarily examined the “demand side” of the student–school relationship, focusing on student participation. These studies typically survey individual students and analyze how factors such as school size (Kleinert, 1969; McNeal, 1999) and socioeconomic status (Feldman & Matjasko, 2007) influence participation rates. In contrast, this study takes a different approach by focusing on the “supply side,” specifically the number of clubs offered by different schools and the impact of various factors (e.g., school enrollment, household income, pupil-to-teacher ratio, and racial demographics) on the club count. This previously overlooked supply-side information is important because, for example, among schools with similar enrollment and socioeconomic status, one offering significantly fewer clubs (i.e., fewer supplies) may experience reduced student participation rates, regardless of students’ potential interest (i.e., potential demand).

Second, previous studies often rely on data collected from only a few schools through in-person surveys, resulting in insufficient data to draw statistically significant conclusions. In contrast, leveraging the public data accessibility in the Internet era, this study adopts a novel approach to collect and clean club data from hundreds of schools’ websites. This dataset, which is two orders of magnitude larger than those in previous studies, enables a more comprehensive and in-depth analysis.

Third, previous studies in social science often assume a simple linear relationship between predictors (e.g., a school’s club count) and independent variables (e.g., school enrollment and socioeconomic status), typically relying on correlation analysis. However, we find that higher-order effects and nonlinearity are present among the factors we study. Empowered by the larger dataset, this study holistically evaluates various linear and nonlinear models to identify the one that most accurately captures the complex relationship among these factors.

In the following sections, we describe how this study addresses gaps identified in the previous research. Before delving into the details, we briefly summarize the key findings for ease of reference:

- (1) The factors influencing the number of clubs offered by schools, ranked by impact, are school size, household income, and pupil-to-teacher ratio.
- (2) After accounting for their indirect influence through school size and household income, racial demographics do not significantly affect club count.
- (3) The number of clubs offered by a school is correlated with  $\sqrt{X_{\text{enroll}}}$ , where  $X_{\text{enroll}}$  is the school’s enrollment.

- (4)  $\sqrt{X_{\text{enroll}}}$  has a multiplicative effect on the influence of a school’s household income and pupil-to-teacher ratio on its club count.
- (5) Schools with a lower pupil-to-teacher ratio tend to offer more clubs.
- (6) Lower-income schools tend to offer fewer clubs, but this effect flattens out when  $X_{\text{lunch}}$  exceeds 40%, where  $X_{\text{lunch}}$  denotes the fraction of students receiving free or reduced-price lunch. For example, while a school with  $X_{\text{lunch}} = 10\%$  tends to offer more clubs than one with  $X_{\text{lunch}} = 40\%$ , assuming other factors are comparable, a school with  $X_{\text{lunch}} = 40\%$  is not expected to offer more clubs than one with  $X_{\text{lunch}} = 80\%$ , since the impact of  $X_{\text{lunch}}$  tends to flatten beyond 40%.
- (7) Among schools with similar demographics (enrollment, household income, and pupil-to-teacher ratio), the top quartile offers 3.2 times more clubs than the bottom quartile, emphasizing the pivotal role of individual school initiatives.

## 2 Methodology for Data Collection

American high schools are secondary schools typically covering grades 9–12. Nearly every high school offers certain school-sponsored extracurricular clubs overseen by the school administration. These clubs provide students with opportunities to engage in a wide variety of activities, including STEM (e.g., robotics), career path (e.g., Future Farmers of America), arts (e.g., drama), social issues (e.g., LGBTQ+ rights), student association (e.g., Black Student Union), leadership and community service (e.g., Key Club), political awareness (e.g., Junior State of America), competition (e.g., Model United Nations), school-community building (e.g., newspaper), culture (e.g., Chinese traditional clothing), intramural sports (e.g., mountain biking), hobbies (e.g., chess), and religions (e.g., Catholicism).

The quantity and variety of clubs provided by high schools are influenced by various factors, ultimately affecting students’ participation. For instance, smaller schools might struggle to gather enough students interested in starting a Latin club, while larger schools might have a surplus of students vying for club leadership positions, resulting in the overmanning effect (Barker & Gump, 1964). Moreover, schools with a high pupil-to-teacher ratio, indicating fewer teachers, may struggle to find a teacher with the necessary bandwidth or expertise to oversee a specific club. Finally, schools predominantly composed of students from lower-income families might face financial

constraints preventing the establishment of clubs like robotics.

We collected a substantial amount of data to analyze the impact of various factors on schools' club counts. Specifically, we used data from the American National Center for Education Statistics (NCES) (NCES, 2023) for information about school sizes, free or reduced-price lunch, pupil-to-teacher ratios, and racial demographics. In addition, we collected complete lists of clubs offered by many schools.

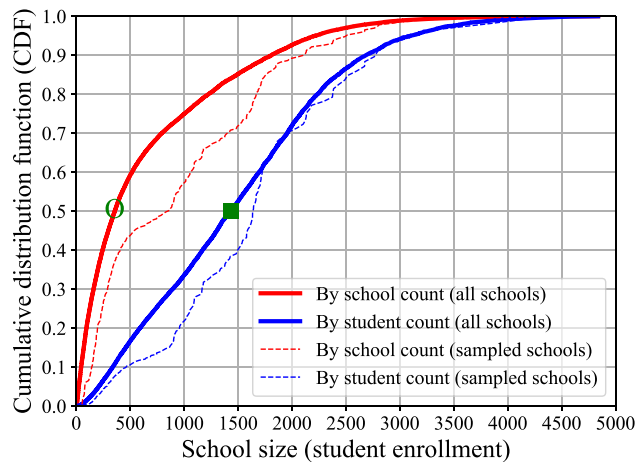
Of all the American public schools that have data reported by NCES for the school year 2021–2022, we selected high schools for this study based on the following criteria. We excluded schools that do not offer grade 12 and only counted students in grades 9–12. Throughout the rest of this article, the term “students” only refers to grade 9–12 students in a school, excluding students from other grades, if present. We focused on physical schools and excluded virtual schools. In addition, we excluded private schools. Finally, we excluded schools that provided no information about their free or reduced-price lunch program, as well as schools that did not report their pupil-to-teacher ratio. This resulted in a total of 21,064 schools that enrolled a total of 14,030,023 grade 9–12 students. We refer to these schools as “candidate schools” or “all schools” throughout the rest of this article.

## 2.1 Sampling of Schools

We first present the characteristics of schools, followed by a description of our school sampling methodology. Figure 1 shows the cumulative distribution of school count and student count as a function of school size. The “all schools” curves represent the 21,064 schools. The “sampled schools” curves represent 229 schools that were randomly selected from the 21,064 schools for this study. We will describe the selection of these schools later.

As an example of how to read the figure, the data point in the green circle on the “by school count (all schools)” curve means that 50% (see the Y-axis) of all schools have 357 (see the X-axis) or fewer students. Similarly, the data point in the green square on the “by student count (all schools)” curve means that 50% of all students across all schools are enrolled in schools that each have 1,434 or fewer students.

These data highlight a significant divergence between school count and student count. While the majority of schools are small, the majority of students are enrolled in larger schools. This affects how we select the sampled schools. Initially, we selected the sampled schools in a way



**Figure 1:** Cumulative distribution of school count and student count as functions of school size. While most schools are small, most students attend larger schools. To avoid a uniform sample dominated by small schools, we intentionally oversampled larger ones.

that ensured, in Figure 1, the curve of “by school count (sampled schools)” closely matched the curve of “by school count (all schools).” However, this led to a lack of sufficient data points for large schools because small schools make up the majority. Specifically, schools with 1,000 or fewer students account for 75% of all schools while hosting only 34% of all students. To address this issue, we intentionally sampled more schools of larger sizes. This leads to the divergence of the “all schools” curves and the “sampled schools” curves in the figure. We will delve further into this topic during the data analysis.

## 2.2 Collecting Club Data

The NCES data contain no information about extracurricular clubs offered by schools. To collect such data, traditional approaches would typically resort to surveys to contact individuals at specific schools, with approvals from numerous ethics committees. However, this approach is challenging to scale to hundreds of schools nationwide, which might be a key reason why, after more than half a century of research on extracurricular activities, there are still no publicly available large-scale datasets providing complete lists of clubs offered by many schools.

In the Internet era, we instead take an innovative approach to collect club data from schools' public websites. Interestingly, the shift to online learning during the COVID-19 pandemic facilitated this data collection process, as an increasing number of schools have adopted websites as a primary channel for communication.

Specifically, we sampled thousands of schools, searched their respective websites, and identified 229 schools with complete lists of clubs they offer. It is important to note that while many school websites include information about a limited subset of clubs as examples, they do not present a comprehensive list of all available clubs. We have excluded these schools from our study. On average, it takes searching through more than 15 schools to identify one with sufficiently complete club data. Consequently, we conducted searches across thousands of schools to identify the 229 sampled schools. In total, the 229 schools hosted 221,300 students and offered 5,983 clubs.

Among a school's clubs, this study excludes clubs for varsity sports such as swimming and baseball for two reasons: (1) They form a substantial category on their own and merit a separate, dedicated study, and (2) the treatment of varsity sports on school websites is highly inconsistent. While some schools classify them as clubs, others categorize them separately under athletics. Furthermore, common varsity sports are frequently not reported on websites despite the likelihood that most schools offer them. For the sake of maintaining consistency in this study, we have excluded varsity sports and intend to explore them in future work. However, we do include intramural sports in this study, such as trapshooting and mountain biking, as they tend to be consistently reported as clubs on different school websites. Similarly, we excluded student councils and class-specific clubs such as Class 2025 because almost every school has these clubs, but the practice of whether to include them as clubs on school websites is inconsistent.

As several co-authors extracted club data from school websites, one challenge was ensuring inter-rater consistency. To address this, we established clear rules to exclude certain schools (e.g., private schools) and adopted a two-step process. First, a software tool randomly sampled thousands of schools from the NCES dataset. Co-authors then searched for each school's website, identified the specific pages listing club data, and excluded schools that listed only a subset of clubs as examples rather than presenting a comprehensive list of all available clubs. This step filtered out approximately 95% of

the sampled schools, leaving around 250. In the second step, a single evaluator reviewed the club webpages of the approximately 250 remaining schools and further excluded any that appeared to have incomplete club data, resulting in a final set of 229 schools. Having a single evaluator performing the second step ensured consistency and avoided inter-rater discrepancies in the final list. The process is scalable, as over 95% of the effort occurred in the first step, which was parallelized across multiple co-authors. A limitation, however, is that some valid schools may have been prematurely excluded in the first step by different data collectors, without a chance for the final evaluator to review them – potentially introducing bias and reducing the dataset size.

### 3 Data Analysis and Findings

In this section, we analyze the effect of various factors on the number of clubs offered by schools. For convenience, the variables used in our analysis are summarized in Table 1, and the basic school statistics are summarized in Table 2.

Our goal is to develop a multiple regression model for predicting the number of clubs ( $Y_{\text{club}}$ ) offered by a school, based on a set of independent variables, including school size ( $X_{\text{enroll}}$ ), the proportion of students receiving free or reduced-price lunch ( $X_{\text{lunch}}$ ), the pupil-teacher ratio ( $X_{\text{teacher}}$ ), and the proportions of students from various racial backgrounds ( $X_{\text{race}}$ , such as  $X_{\text{white}}$ ,  $X_{\text{hispanic}}$ ,  $X_{\text{black}}$ , etc.). Rather than blindly dumping all these variables into a regression model, we first thoroughly examine the characteristics of each variable to guide the design of the regression model. This helps us discover and address potential issues such as nonlinearity, heteroscedasticity, and multicollinearity.

#### 3.1 Effect of School Size

Our analysis starts with the school size factor. For an initial intuitive grasp of the effect of school size, we plot in Figure 2

**Table 1:** Symbols for dependent and independent variables used in this study

Variable	Explanation
$Y_{\text{club}}$	Number of clubs offered by a school
$X_{\text{enroll}}$	School size, i.e., student enrollment
$X_{\text{teacher}}$	Pupil-to-teacher ratio. A higher $X_{\text{teacher}}$ means fewer teachers per student
$X_{\text{lunch}}$	Fraction of students in a school receiving free or reduced-price lunch. A higher $X_{\text{lunch}}$ value indicates a lower household income
$X_{\text{race}}$ ( $X_{\text{white}}$ , $X_{\text{hispanic}}$ , $X_{\text{black}}$ , etc.)	Fraction of students of a specific race in a school

**Table 2:** Average school statistics. For example, the “White” column shows that, among all 21,064 schools, 45.45% of students are White, while among the 229 sampled schools, 41.57% of students are White. This table shows that the sampled schools are representative in terms of racial demographics

Variable	$X_{\text{enroll}}$	$X_{\text{lunch}}$ (%)	$X_{\text{teacher}}$	White (%)	Hispanic (%)	Black (%)	Asian/Pacific Islander (%)	Multi-race (%)	Native American (%)	Hawaiian or other Pacific Isl. (%)
All schools	666	50	15.0	45.45	28.80	15.03	5.48	3.92	0.88	0.38
Sampled schools	966	47	15.3	41.57	31.36	14.10	7.91	4.02	0.56	0.41

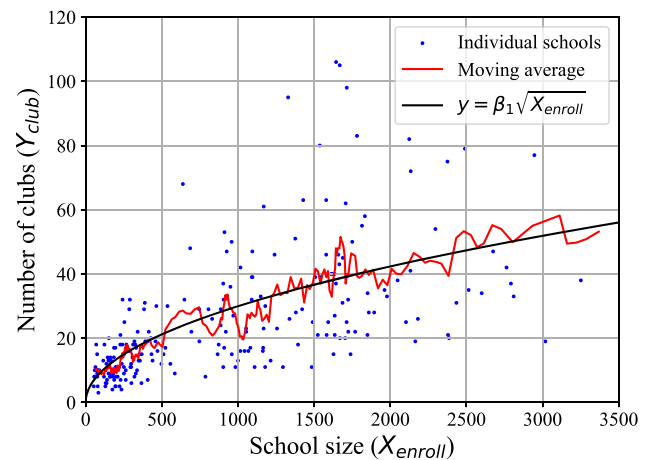
the relationship between school size ( $X_{\text{enroll}}$ ) and the number of clubs ( $Y_{\text{club}}$ ) for each of the 229 sampled schools. Each dot represents one school, and the “moving average” curve shows the overall trend. To compute the moving-average curve, we first sort the schools according to their sizes. Then, for each group of 10 consecutive schools in the sorted list, we calculate their average school size ( $x$ ) and average number of clubs ( $y$ ). These mean values are then plotted as one data point ( $x, y$ ) on the moving-average curve. By iterating this procedure for all groups of 10 consecutive schools in the sorted list, the complete moving-average curve can be plotted.

The moving-average curve reveals a strong positive correlation between  $X_{\text{enroll}}$  and  $Y_{\text{club}}$ . The Pearson correlation coefficient is 0.63, with  $p < 0.001$ . To further understand their relationship, we fit the data to the following model:  $Y_{\text{club}} = \beta_1 f(X_{\text{enroll}}) + \varepsilon$ . Note that this model does not have a bias term because, intuitively, when  $X_{\text{enroll}}$  approaches zero,  $Y_{\text{club}}$  should also approach zero. This is confirmed by the trend in Figure 2.

The moving-average curve in Figure 2 shows that  $Y_{\text{club}}$  grows more slowly than  $X_{\text{enroll}}$ . This pattern suggests modeling the relationship with a slow-growing function. After exploring various forms of function  $f(\cdot)$ , including  $f(x) = x$ ,  $f(x) = \sqrt{x}$ ,  $f(x) = \sqrt[3]{x}$ ,  $f(x) = \log(x)$ , and a logistic function, we find that  $f(x) = \sqrt{x}$  yields the best fit, specifically, with the following coefficient:

$$Y_{\text{club}} = \beta_1 f(X_{\text{enroll}}) + \varepsilon = 0.95\sqrt{X_{\text{enroll}}} + \varepsilon \approx \sqrt{X_{\text{enroll}}}. \quad (1)$$

This function is also plotted in Figure 2, showing good alignment with the moving-average curve. Note that we



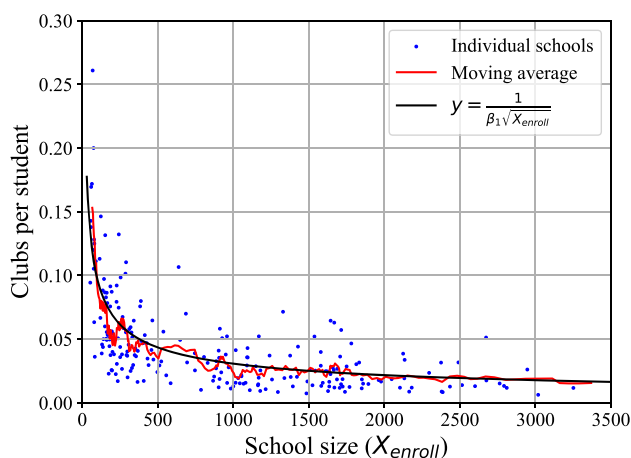
**Figure 2:** The number of clubs follows the trend of square root of school size. Without thorough modeling and exploration of different curve-fitting functions, one might oversimplify the relationship between school size and number of clubs, mistakenly assuming a linear correlation that could lead to suboptimal results.

fit the function to the entire dataset instead of the moving-average curve, which merely aids visualization.

We derive several observations from Figure 2. First, the simplicity of the function is appealing, indicating a square root relationship between  $Y_{club}$  and  $X_{enroll}$ . Second, it reveals the nonlinear nature of their relationship, highlighting the necessity for a square root transformation of  $X_{enroll}$  in the final multiple regression model. Finally, as the school size increases, the prediction error also increases, indicating the presence of heteroscedasticity. This issue requires attention in the final multiple regression model.

In Figure 2, one might note that a big fraction of the sampled schools are relatively small, i.e., with 500 or fewer students. This is because most schools are smaller schools, as shown by the curve of “by school count (all schools)” in Figure 1. To have sufficient samples for larger schools, we already intentionally sampled a bigger fraction of larger schools, as shown by the curve of “by school count (sampled schools)” in Figure 1.

Finally, from Figure 2, one might be inclined to assume that students at smaller schools are at a disadvantage, given their schools offer fewer clubs. However, Figure 3 presents a contrasting perspective, unveiling that the number of clubs per student is, in fact, higher in smaller schools. Because the number of clubs follows the trend of  $\sqrt{X_{enroll}}$ , the number of clubs per student follows the trend of  $\frac{\sqrt{X_{enroll}}}{X_{enroll}} = \frac{1}{\sqrt{X_{enroll}}}$ . On the one hand, students at smaller schools enjoy better opportunities for leadership roles and active club engagement. On the other hand, smaller schools have a reduced variety of club choices.



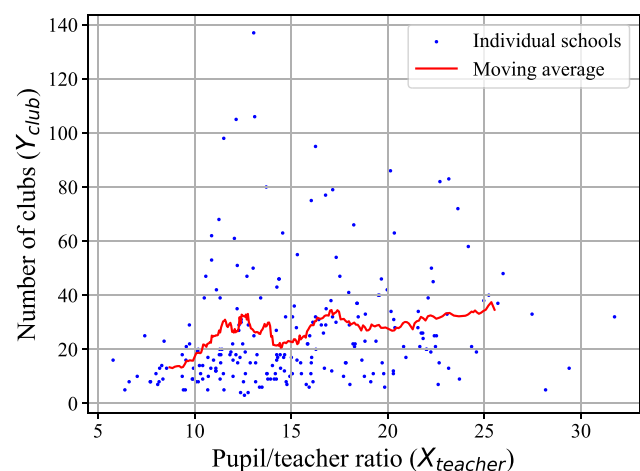
**Figure 3:** The number of clubs per student decreases as the school size increases. On the one hand, smaller schools tend to have fewer clubs, limiting students' choices. On the other hand, students at smaller schools enjoy better opportunities for leadership roles and active club engagement.

### 3.2 Effect of Pupil-to-Teacher Ratio

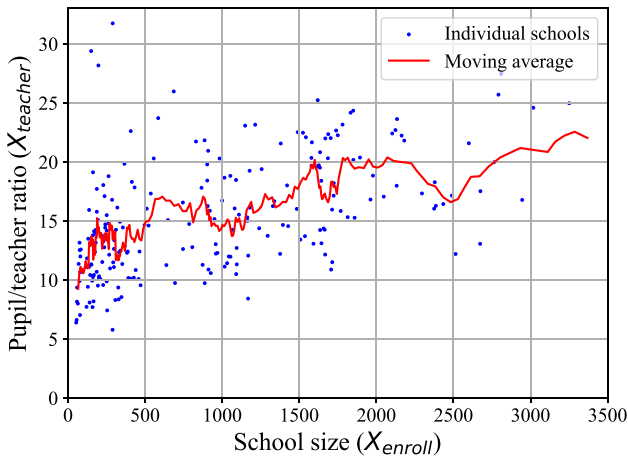
In addition to school size, the pupil-to-teacher ratio ( $X_{teacher}$ ) also affects the number of clubs. For an initial intuitive grasp, we graph in Figure 4 the 229 sampled schools based on their  $X_{teacher}$  and  $X_{enroll}$  values. Surprisingly, the figure appears to imply a positive correlation between  $X_{teacher}$  and  $Y_{club}$ , which is confirmed by a correlation coefficient of 0.20 and  $p < 0.01$ . It seems counterintuitive for schools with a higher  $X_{teacher}$  (meaning fewer teachers per student) to have more clubs, as schools rely on teachers to assume advisory roles for each club.

Further investigation reveals that this counterintuitive result is merely a side effect of the correlation between  $X_{teacher}$  and  $X_{enroll}$ , as depicted in Figure 5. The seemingly paradoxical relationship between  $X_{teacher}$  and  $Y_{club}$  is primarily driven by the fact that larger schools, with higher  $X_{teacher}$ , tend to have more clubs, rather than the direct effect of  $X_{teacher}$  on  $Y_{club}$ .

To assess the true effect of  $X_{teacher}$ , we design a  $t$ -test between two groups of schools categorized by higher and lower  $X_{teacher}$  values, respectively, while ensuring the two groups have comparable characteristics in school size. As portrayed in Figure 5, smaller schools generally exhibit lower  $X_{teacher}$  values. Consequently, if we were to partition the groups solely based on raw  $X_{teacher}$  values, the lower- $X_{teacher}$  group would largely consist of smaller schools, and the higher- $X_{teacher}$  group would mainly include larger schools. This would result in a comparison predominantly between schools of differing sizes, which is undesirable.



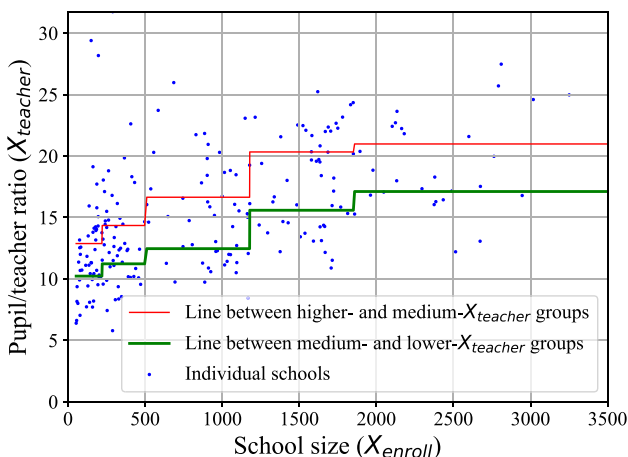
**Figure 4:** Schools with higher  $X_{teacher}$  values tend to have more clubs, which is counterintuitive at first glance. This is because larger schools tend to be associated with both higher  $X_{teacher}$  and higher  $Y_{club}$  values, leading to the side effect of a positive correlation between  $X_{teacher}$  and  $Y_{club}$ .



**Figure 5:**  $X_{teacher}$  is correlated with  $X_{enroll}$ . Larger schools tend to have higher  $X_{teacher}$  values.

To prevent this, we follow the division lines in Figure 6 to create school groups for comparison. Specifically, we first sort the schools by their sizes. For every 30 adjacent schools in the sorted list with similar school sizes, we identify division lines that divide the 30 schools into three groups of equal sizes with higher, medium, and lower  $X_{teacher}$  values, respectively. As observed in the figure, these division lines ascend as the school size increases, accommodating the tendency for larger schools to have higher  $X_{teacher}$  values. Through this approach, schools assigned to the higher- and lower- $X_{teacher}$  groups possess comparable  $X_{enroll}$  values, averaging 790 and 879, respectively. In contrast, their  $X_{teacher}$  values differ significantly, averaging 19 and 11, respectively. These characteristics align well with the experiment's objectives.

We compare the higher- and lower- $X_{teacher}$  groups in Figure 7. The figure shows that, after excluding the



**Figure 6:** Partition the 229 sampled schools into three groups based on their  $X_{teacher}$  values.

influence of school size, schools in the lower- $X_{teacher}$  group tend to have more clubs. A  $t$ -test between the two groups validates that the difference is statistically significant, with  $t$ -statistic =  $-2.42$  and  $p < 0.02$ . On average, schools in the lower- $X_{teacher}$  group have 9.3 more clubs than those in the higher- $X_{teacher}$  group. This contrasts with the seemingly positive correlation between  $X_{teacher}$  and  $Y_{club}$ .

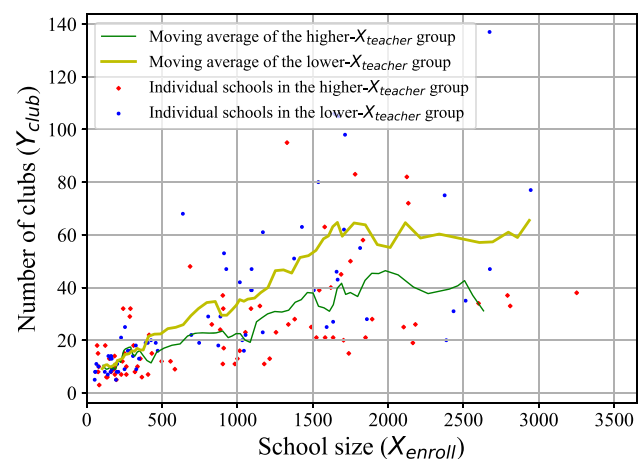
Furthermore, Figure 7 illustrates that the gap in  $Y_{club}$  values between the two groups tends to widen as  $X_{enroll}$  increases. This suggests that  $X_{enroll}$  magnifies the effect of  $X_{teacher}$ , which is why we will employ  $X_{enroll}$  as a multiplicative factor for other independent variables such as  $X_{teacher}$  in the final multiple regression model.

### 3.3 Effect of Household Income

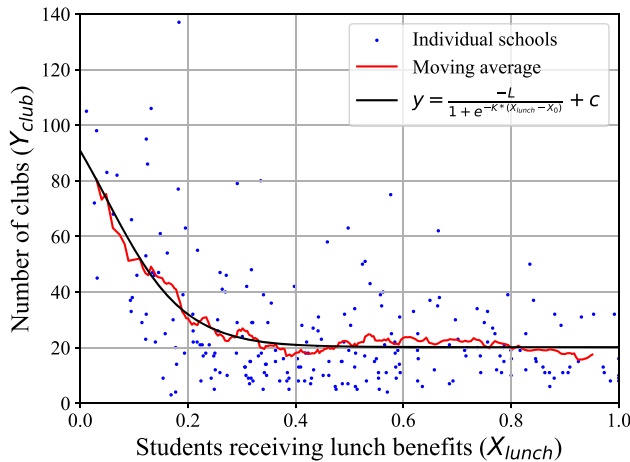
Next, we analyze the effect of household income by using  $X_{lunch}$  as an indicator of household income.  $X_{lunch}$  is calculated as  $l/X_{enroll}$ , where  $l$  is the number of students eligible for free or reduced-price lunch. A lower  $X_{lunch}$  indicates a higher income.

To assess whether we need to exclude the influence of school size before analyzing the effect of  $X_{lunch}$ , we calculate the correlation between  $X_{lunch}$  and school size. Unlike school size's strong correlation with  $X_{teacher}$ , it has little correlation with  $X_{lunch}$ , reflected in a correlation coefficient of  $-0.08$  and  $p = 0.22$ . This suggests that we can directly observe the effect of  $X_{lunch}$  on  $Y_{club}$  without needing to exclude the influence of school size.

Figure 8 illustrates the relationship between  $X_{lunch}$  and  $Y_{club}$ . When  $X_{lunch}$  is less than 40%, an increase in  $X_{lunch}$  tends to cause a decrease in  $Y_{club}$ . However, when  $X_{lunch}$



**Figure 7:** Schools in the lower- $X_{teacher}$  group tend to have more clubs, meaning that among schools of similar sizes, those with a lower pupil-to-teacher ratio tend to have more clubs.



**Figure 8:** Only when  $X_{\text{lunch}} < 40\%$ , higher-income schools tend to have more clubs than lower-income schools.

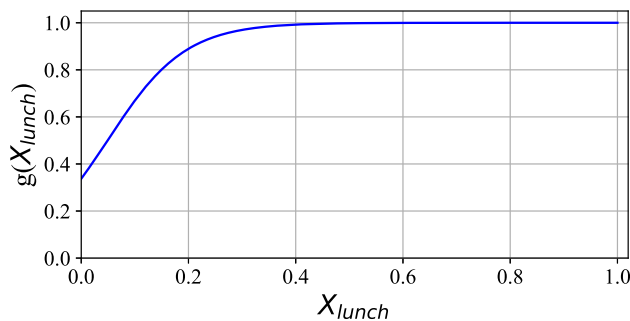
exceeds 40%, a further increase in  $X_{\text{lunch}}$  no longer affects  $Y_{\text{club}}$ . In general, the logistic function – visualized in Figure 9 and defined mathematically in equation (9) – is commonly used to represent this flattening-out effect. Specifically, the logistic function depicted in Figure 8 aligns well with the moving-average curve. Its parameters,  $K$ ,  $X_0$ ,  $L$ , and  $C$ , are inferred using the Levenberg–Marquardt algorithm to fit the data. Our empirical evaluation further demonstrates that among various functions we explored, the logistic function indeed offers the best fit.

Because of the nonlinear relationship between  $X_{\text{lunch}}$  and  $Y_{\text{club}}$ , we will apply the following nonlinear transformation to  $X_{\text{lunch}}$  in the final multiple regression model.

$$g(X_{\text{lunch}}) = \frac{1}{1 + e^{-K(X_{\text{lunch}} - X_0)}} = \frac{1}{1 + e^{-13.8(X_{\text{lunch}} - 0.049)}}. \quad (2)$$

$g(X_{\text{lunch}})$  is graphed in Figure 9. Within the domain of  $[0, 1]$  for  $X_{\text{lunch}}$ , the range of  $g(X_{\text{lunch}})$  starts at 0.337 and levels off at 1 when  $X_{\text{lunch}}$  approaches approximately 0.4.

Overall, because the curve in Figure 8 flattens out when  $X_{\text{lunch}}$  exceeds 40%, it suggests that in terms of the



**Figure 9:** This logistic function  $g(X_{\text{lunch}})$  is used to transform  $X_{\text{lunch}}$  in the final multiple regression model.

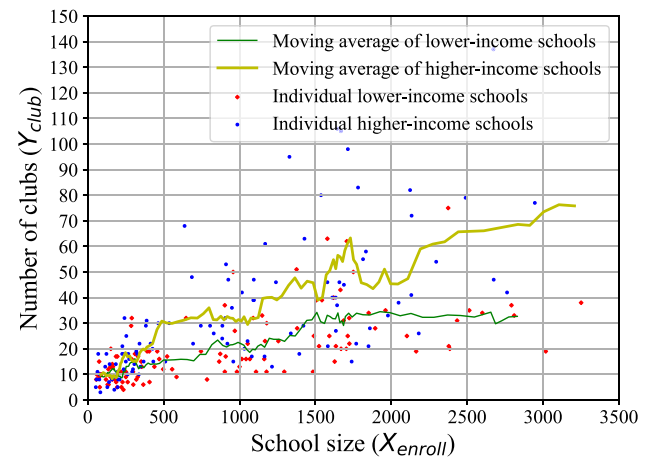
effect on club count, there is a significant difference between the higher- and middle-income groups, but there is little difference between the middle- and lower-income groups. For context, the average value of  $X_{\text{lunch}}$  for all schools is 50%.

To further quantify the statistical significance of the effect of household income, we follow the method shown in Figure 6 to partition the 229 sampled schools into three groups based on the values of  $X_{\text{lunch}}$ : one-third higher-income schools (i.e., lower  $X_{\text{lunch}}$ ), one-third middle-income schools, and one-third lower-income schools. We compare  $Y_{\text{club}}$  of the higher- and lower-income groups, as depicted in Figure 10. The figure shows that the higher-income group has more clubs than the lower-income group, specifically, with an average difference of 16.7 clubs. For context, on average, a school has 26.1 clubs. A  $t$ -test between these groups confirms that their difference is statistically significant, with  $t$ -statistics =  $-4.77$  and  $p < 0.0001$ .

Similar to Figures 7 and 10 also shows that the disparity in  $Y_{\text{club}}$  values between the two groups tends to expand as  $X_{\text{enroll}}$  increases. This implies that  $X_{\text{enroll}}$  amplifies the influence of  $X_{\text{lunch}}$ , which is why we will include  $X_{\text{enroll}}$  as a multiplicative factor for other independent variables in the final multiple regression model.

### 3.4 Effect of Racial Demographics

In this section, we analyze the effect of racial demographics. We begin by examining the largest racial group, the White students. For each school, we compute the fraction of White students as  $X_{\text{white}} = w/X_{\text{enroll}}$ , where  $w$  is the number of White students in the school.  $X_{\text{white}}$  is correlated



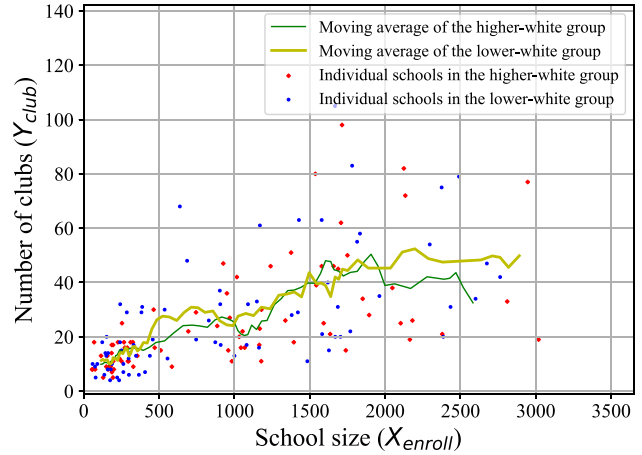
**Figure 10:** Higher-income schools tend to have more clubs.

with school size, with a coefficient of  $-0.34$  and  $p < 0.0001$ . This indicates that larger schools tend to have smaller proportions of White students. Moreover,  $X_{\text{white}}$  is even more strongly correlated with  $X_{\text{lunch}}$ , with a coefficient of  $-0.54$  and  $p < 0.0001$ . This indicates that lower-income schools tend to have smaller proportions of White students.

Since both  $X_{\text{enroll}}$  and  $X_{\text{lunch}}$  influence how  $X_{\text{white}}$  affects  $Y_{\text{club}}$ , to assess the true effect of  $X_{\text{white}}$  in a controlled environment, we design an experiment to compare two groups of schools with comparable values in both  $X_{\text{lunch}}$  and  $X_{\text{enroll}}$ , but differ significantly in their  $X_{\text{white}}$  values. The subsequent procedure outlines our approach. First, we sort the 229 sampled schools by their  $X_{\text{lunch}}$  values and partition the sorted list into  $N_{\text{lunch}}$  buckets. Each bucket comprises the same number of schools, and the schools within the same bucket have similar  $X_{\text{lunch}}$  values as they are adjacent on the list sorted by  $X_{\text{lunch}}$ . For schools in each bucket, we resort them by their  $X_{\text{enroll}}$  values and further partition them into  $N_{\text{enroll}}$  bins. Each bin comprises the same number of schools, and the schools within the same bin have similar  $X_{\text{enroll}}$  values as they are adjacent on the list sorted by  $X_{\text{enroll}}$ . In total, there are  $N_{\text{lunch}} \times N_{\text{enroll}}$  bins, where each bin has the same number of schools, and within each bin, schools exhibit similar values in both  $X_{\text{lunch}}$  and  $X_{\text{enroll}}$ . Finally, for the schools in each bin, we partition them into three subgroups based on their  $X_{\text{white}}$  values, forming the higher-, medium-, and lower- $X_{\text{white}}$  subgroups. The aggregation of all schools in the higher- $X_{\text{white}}$  subgroups across all  $N_{\text{lunch}} \times N_{\text{enroll}}$  bins forms the overall higher- $X_{\text{white}}$  group. Similarly, the assembly of all schools in the  $N_{\text{lunch}} \times N_{\text{enroll}}$  lower- $X_{\text{white}}$  subgroups forms the overall lower- $X_{\text{white}}$  group. In our experiment, both  $N_{\text{lunch}}$  and  $N_{\text{enroll}}$  are set to 5.

The resulting higher- and lower- $X_{\text{white}}$  groups exhibit desired characteristics. Schools in these two groups have comparable  $X_{\text{enroll}}$  values, averaging 969 and 994 students, respectively, along with comparable  $X_{\text{lunch}}$  values, averaging 46 and 49%, respectively. However, their  $X_{\text{white}}$  values differ significantly at 70 and 28%, respectively – meeting our intended criteria. Remarkably, despite substantial  $X_{\text{white}}$  discrepancies, they have a comparable number of clubs ( $Y_{\text{club}}$ ), averaging 25.4 and 27.5, respectively. Interestingly, the higher- $X_{\text{white}}$  group shows a slightly lower number of clubs. The  $t$ -test between the  $Y_{\text{club}}$  values of these two groups yields a  $p$ -value of 0.19, insufficient to establish a significant difference between them.

To visualize the pattern, we plot the schools in the higher- and lower- $X_{\text{white}}$  groups in Figure 11. The figure shows that despite the notable difference in their  $X_{\text{white}}$  values, there is no distinct disparity in  $Y_{\text{club}}$  between the two groups.



**Figure 11:** For two school groups with comparable  $X_{\text{enroll}}$  and  $X_{\text{lunch}}$  values but significant differences in their  $X_{\text{white}}$  values, there is no evidence that  $X_{\text{white}}$  significantly affects  $Y_{\text{club}}$ .

Similar  $t$ -test results for other racial groups are summarized in Table 3. Overall, when comparing two groups of schools with higher and lower fractions of students from a specific racial group, similar to what is depicted in Figure 11, the marginal difference in  $Y_{\text{club}}$  and the large  $p$  value in the  $t$ -test indicate a lack of statistical significance in the difference. An intuitive interpretation of these results is that schools with comparable student enrollments and household incomes do not display significant disparities in the number of clubs due to variations in racial demographics. The lack of direct affect on  $Y_{\text{club}}$  is a key reason for us to consider excluding the factors representing racial demographics from the final multiple regression model, as described in the next section.

### 3.5 Multiple Regression Model

After understanding the characteristics of individual independent variables, in this section, we design a multiple regression model that incorporates all the factors. Specifically, we model the number of clubs as

$$Y_{\text{club}} = f(X_{\text{enroll}}) \left( \beta_0 + \beta_1 \cdot g(X_{\text{lunch}}) + \beta_2 \cdot X_{\text{teacher}} + \sum_{\text{race}} \beta_{\text{race}} \cdot X_{\text{race}} \right) + \varepsilon, \quad (3)$$

where  $\beta_i$  are coefficients,  $\varepsilon$  is the model residual,  $X_{\text{race}}$  such as  $X_{\text{white}}$  is the fraction of students of a specific race in a school,  $f(X_{\text{enroll}})$  is the nonlinear transformation of  $X_{\text{enroll}}$  derived from equation (1), and  $g(X_{\text{lunch}})$  is the nonlinear transformation of  $X_{\text{lunch}}$  obtained from equation (2).

**Table 3:** Results of  $t$ -tests that explicitly compare two groups of schools with higher and lower values on a specific independent variable

	$X_{\text{teacher}}$	$X_{\text{lunch}}$	White ( $X_{\text{white}}$ )	Hispanic	Black	Asian/Pacific Islander	Multirace	Native American	Hawaiian/other Pacific Isl.
Difference in mean	-9.3	-16.7	-2.1	5.4	-1.8	0.1	0.0	-4.7	-1.6
$t$ -statistics	-2.42	-4.77	-0.64	1.59	-0.51	0.03	-0.01	-1.31	-0.49
$p$ -value	< 0.02	<0.0001	0.5203	0.114	0.6103	0.9792	0.9955	0.1922	0.6271

The  $X_{\text{teacher}}$  column corresponds to the comparison shown in Figure 7. The  $X_{\text{lunch}}$  column corresponds to the comparison shown in Figure 10. The  $X_{\text{white}}$  column corresponds to the comparison shown in Figure 11. The other  $X_{\text{race}}$  columns are obtained through experiments similar to the one shown in Figure 11.

$$\begin{aligned} f(X_{\text{enroll}}) &= \sqrt{X_{\text{enroll}}} \\ g(X_{\text{lunch}}) &= \frac{1}{1 + e^{-13.8(X_{\text{lunch}} - 0.049)}}. \end{aligned} \quad (4)$$

A distinctive aspect of this model is the role of  $f(X_{\text{enroll}})$  as a multiplicative factor for all other independent variables. This choice is closely tied to the significant role of school size. When  $f(X_{\text{enroll}})$  approaches 0, indicating a small student population,  $Y_{\text{club}}$  should converge to zero, regardless of  $X_{\text{teacher}}$ ,  $X_{\text{lunch}}$ , and  $X_{\text{race}}$ . This behavior is effectively captured by utilizing  $f(X_{\text{enroll}})$  as a multiplicative term. Furthermore, the effect of other independent variables is amplified by school size. For instance, if higher income levels are assumed to lead to more clubs, this effect should result in a more substantial absolute increase in the number of clubs in larger schools. This amplification finds empirical support in the data, as depicted in Figures 7 and 10. Also note that this model does not have a bias term in order to achieve the effect that as  $f(X_{\text{enroll}})$  approaches 0,  $Y_{\text{club}}$  converges to zero.

In addition to this model, we have explored numerous others, including those incorporating additional terms related to  $X_{\text{lunch}}$  and  $X_{\text{teacher}}$  but without  $f(X_{\text{enroll}})$  as their multiplicative factor, or adding a bias term, etc. Moreover, we have also explored introducing independent variables besides school demographics, such as those representing how schools present their club information on their websites, with or without teacher contact information, and

with or without detailed club descriptions. Nevertheless, these models added complexity without meaningfully improving data fitting. Consequently, we have excluded their use.

### 3.6 Further Simplifying the Model

Although the model in equation (3) seems intuitive and comprehensive, we find that the terms representing racial demographics,  $\sum_{\text{race}} \beta_{\text{race}} X_{\text{race}}$ , do not contribute significantly to enhancing data fitting. Therefore, we eliminate those terms and adopt the simplified model below:

$$Y_{\text{club}} = f(X_{\text{enroll}}) \cdot (\beta_0 + \beta_1 \cdot g(X_{\text{lunch}}) + \beta_2 \cdot X_{\text{teacher}}) + \varepsilon. \quad (5)$$

We use the linear regression implementation in the statsmodels package (Statsmodels, 2023) to compute the coefficients. Information concerning these coefficients is summarized in Table 4. Due to the presence of heteroscedasticity in the residual, as evident in Figures 7 and 10, we employ the HC3 covariance matrix estimator in statsmodels to calculate heteroskedasticity-robust standard errors. This estimator implements the algorithm proposed by MacKinnon and White (MacKinnon & White, 1985).

The adjusted  $R^2$  for the simplified model is 0.60, while the  $p$ -value for the F-statistic is less than  $10^{-60}$ , indicating that the model is a reasonable fit for the data. Note that

**Table 4:** Information about the model coefficients in equation (5).  $\beta_2$  is much smaller than  $\beta_1$  due to the considerably larger scale of  $X_{\text{teacher}}$  (mean: 15.3) in comparison to that of  $X_{\text{lunch}}$  (mean: 0.5)

Coefficient symbol	Coefficient value	Robust standard error	$t$ -value	$p$ -value	Confidence interval		Partial $R^2$
					[0.025	0.975]	
$\beta_0$ for $X_{\text{enroll}}$	3.2124	0.333	9.66	<0.001	2.557	3.868	0.50
$\beta_1$ for $g(X_{\text{lunch}})$	-1.9507	0.298	-6.535	<0.001	-2.539	-1.363	0.30
$\beta_2$ for $X_{\text{teacher}}$	-0.0273	0.008	-3.563	<0.001	-0.042	-0.012	0.08

when dealing with models lacking a bias term, the statsmodels package defaults to presenting *uncentered* and adjusted  $R^2$  values, which is 0.84 for our model. When computing  $R^2$ , it calculates the sum of squares total as  $\sum Y_{\text{club}}^2$ , instead of the conventional  $\sum (Y_{\text{club}} - \bar{Y}_{\text{club}})^2$  due to the absence of the bias term to help center the residual. To avoid undue optimism, we have opted to focus on the lower, *centered*  $R^2$  value of 0.60, instead of the *uncentered*  $R^2$  value of 0.84.

The mean of the model residual,  $\varepsilon$ , is only  $-0.44$ , in comparison to the mean of  $Y_{\text{club}}$  at 26.1. This suggests that the decision to exclude a bias term is justified. Furthermore, empirically, we note that the absence of the bias term enhances the model's stability by preventing the regression from adjusting the bias term to compensate for errors elsewhere in the coefficients.

The signs of the coefficients in Table 4 align with the previous analysis results for individual variables. Specifically,  $X_{\text{enroll}}$  shows a positive correlation with  $Y_{\text{club}}$ , while  $X_{\text{lunch}}$  and  $X_{\text{teacher}}$  display negative correlations with  $Y_{\text{club}}$  after excluding the influence of  $X_{\text{enroll}}$ . Furthermore, the partial  $R^2$  values for  $X_{\text{enroll}}$ ,  $X_{\text{lunch}}$ , and  $X_{\text{teacher}}$  are 0.50, 0.30, and 0.08, respectively. These values reflect the relative magnitude of their affect on  $Y_{\text{club}}$ , aligning with our prior evaluation of these variables individually.

After establishing the simplified model in equation (5) as a solid baseline, we delve into the reasoning behind excluding the  $\sum_{\text{race}} \beta_{\text{race}} X_{\text{race}}$  terms. When solely introducing the  $X_{\text{white}}$  variable, without incorporating other  $X_{\text{race}}$  variables, to the simplified model, the adjusted  $R^2$  remains 0.60. This indicates that  $X_{\text{white}}$  does not contribute to enhancing data fit. On the other hand, if all  $X_{\text{race}}$  variables are added to the model, they cause multicollinearity as the  $X_{\text{race}}$  variables always sum up to one and have strong correlations. Hence, at least one  $X_{\text{race}}$  variable should be dropped. Eliminating an  $X_{\text{race}}$  variable representing a racial group with a small school population inadequately mitigates multicollinearity, as other high-value  $X_{\text{race}}$  variables still have strong correlations. The most effective resolution is removing the  $X_{\text{race}}$  variable with the largest value, namely,  $X_{\text{white}}$ . With  $X_{\text{white}}$  eliminated and multicollinearity resolved, incorporating additional  $X_{\text{race}}$  variables, regardless of their combination, never raises the adjusted  $R^2$  beyond 0.61, which is hardly an improvement over the simplified model.

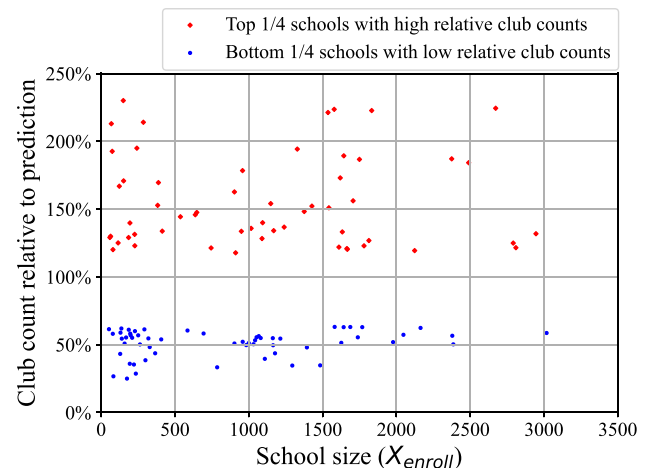
These observations indicate that the inclusion of  $X_{\text{race}}$  variables complicates the model without improving its data fit. Furthermore, the  $t$ -test results in Table 3 explicitly confirm that these variables do not have a significant effect on  $Y_{\text{club}}$  for schools with comparable  $X_{\text{enroll}}$  and  $X_{\text{lunch}}$

values despite the significant variation in the  $X_{\text{race}}$  values. Guided by these insights, we choose to exclude the  $X_{\text{race}}$  terms from the simplified model.

### 3.7 School Initiatives in Improving Club Offerings

On the one hand, the multiple regression model's reasonable accuracy, with an adjusted  $R^2$  of 0.60, represents a noteworthy achievement, considering the intricate interplay of numerous big and small factors affecting operations of schools ranging from 50 to 4,300 students in this study. On the other hand, it underscores that schools still have ample opportunities to take initiatives to enhance club offerings, since 40% of the variance remains unattributed to hard-to-change factors like school size.

To illustrate the substantial effect of school initiatives, we compare top- and bottom-performing schools in terms of club offerings in Figure 12. For each school, we calculate the ratio between its actual club count and the club count prediction derived from equation (5). A school with a higher ratio outperforms a school with a lower ratio. In Figure 12, we plot the top and bottom 1/4 schools with the highest and lowest values in this ratio, respectively. This figure reveals a significant disparity between the top and bottom schools. Specifically, the average of the ratios for the top 1/4 schools is 163%, whereas that for the bottom 1/4



**Figure 12:** Out of the 229 sampled schools, this figure compares the relative club counts between the top 58 schools and the bottom 58 schools. After accounting for school demographics, the top and bottom schools still show a significant difference in club count, emphasizing the pivotal role of individual school initiatives.

schools is only 51%. In other words, the top schools have 3.2 times more clubs than the bottom schools after school demographics are taken into account, emphasizing the pivotal role of individual school initiatives.

### 3.8 Analysis Summary and Lessons Learned

Overall, the simplified model demonstrates a reasonable level of accuracy, indicated by its adjusted  $R^2$  value of 0.60. Notable features of this model include: (1) addressing the nonlinearity of  $X_{\text{enroll}}$  and  $X_{\text{lunch}}$  through the transformations  $f(X_{\text{enroll}})$  and  $g(X_{\text{lunch}})$  as defined in equation (4); (2) utilizing  $X_{\text{enroll}}$  exclusively as a multiplicative factor for  $X_{\text{lunch}}$  and  $X_{\text{teacher}}$ ; (3) refraining from incorporating a bias term to force  $Y_{\text{club}}$  to approach zero as  $X_{\text{enroll}}$  approaches zero; and (4) excluding the  $X_{\text{race}}$  variables from the model despite the common perception that racial demographics might correlate with  $X_{\text{enroll}}$  and  $X_{\text{lunch}}$ , potentially influencing  $Y_{\text{club}}$ .

A valuable lesson we have learned is that the design of this model greatly benefited from a thorough initial analysis of the characteristics of individual variables, rather than blindly including them in a linear regression model. This aided us in identifying and designing the aforementioned features of the model.

In contrast, if the traditional linear regression model shown below were used, it would produce misleading and unreasonable results despite its seemingly acceptable  $R^2$  value of 0.51.

$$Y_{\text{club}} = \beta_0 + \beta_1 X_{\text{enroll}} + \beta_2 X_{\text{lunch}} + \beta_3 X_{\text{teacher}} + \varepsilon. \quad (6)$$

Specifically, the regression result for the aforementioned model misleadingly suggests that  $X_{\text{teacher}}$  may not affect  $Y_{\text{club}}$  due to  $\beta_3$ 's large  $p$ -value of 0.12. More importantly, the bias term,  $\beta_0$ , has a value of 30.0, with  $p < 0.001$ , indicating high confidence. However, this large bias is highly unreasonable since the mean of  $Y_{\text{club}}$  is only 29.1, and when  $X_{\text{enroll}}$  approaches zero,  $Y_{\text{club}}$  should also approach zero, rather than taking on the large bias value of 30.0. These results highlight the importance of introducing the aforementioned notable features in our model.

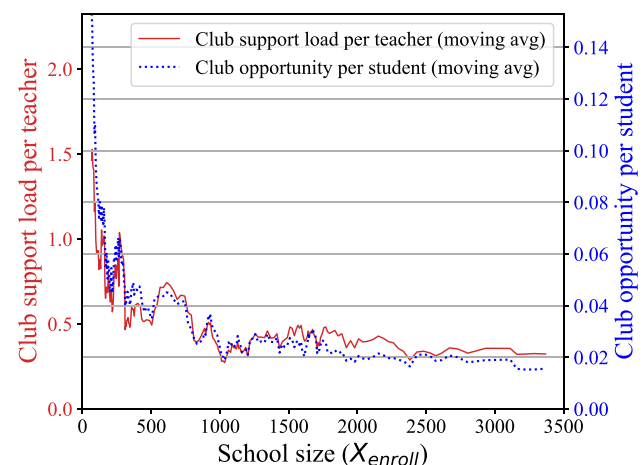
Another valuable lesson we have learned is that when validating observations, we have frequently strived to design experiments that directly compare two population groups differing significantly in one independent variable while maintaining similar characteristics across other independent variables. Concrete examples include the comparisons shown in Figures 7, 10, and 11, which are further summarized in Table 3. This direct validation

approach has enhanced our confidence in and improved our interpretation of the summaries generated by statistical software.

## 4 Discussion and Recommendation

In this section, we discuss issues that may adversely affect the availability of clubs provided by schools, starting with the school size factor, as it has the biggest affect due to its multiplicative role in equation (3). The challenge is evident in Figure 3, which shows a decrease in club opportunity per student as school size increases. In particular, large schools experience the over-manning effect (Barker & Gump, 1964), resulting in a surplus of students competing for club leadership positions, which limits opportunities for leadership skill development. Moreover, as club growth lags behind student enrollment, increased student-to-club ratios discourage participation. These observations align with the prior research indicating lower student engagement in various activities within larger schools (Schoggen & Schoggen, 1988, Morgan & Alwin, 1980, Stevens & Peltier, 1994).

What limits club growth in larger schools? One hypothesis is that, since the pupil-to-teacher ratio tends to be higher in larger schools, as illustrated in Figure 5, the availability of teachers as supervisors for clubs might be a limiting factor. However, Figure 13 contradicts this. As the club opportunity per student decreases with school size, the club support load per teacher similarly decreases, implying teachers are not the bottleneck. In this figure, the moving average curves are computed similarly to those in Figure 2,



**Figure 13:** Club support load per teacher decreases as school size increases.

but individual schools are not plotted to avoid overcrowding the figure. Club support load per teacher is calculated as  $\frac{Y_{\text{club}} \times X_{\text{teacher}}}{X_{\text{enroll}}}$ .

We believe that the inclination to eliminate redundant clubs is a contributing factor that negatively affects club offerings, leading to the emergence of large clubs and the subsequent over-manning effect, ultimately discouraging participation. Comparing club offerings to course offerings, we observe that in large schools, popular courses such as algebra often offer multiple parallel classes to reduce class size and promote engagement. However, this approach is seldom applied to clubs, even though clubs like chess and Key Club are popular. Despite the theoretical possibility of dividing large clubs into smaller subgroups for internal operations, this is rarely practiced. In addition, the number of club officers rarely increases in proportion to club size. These factors naturally restrict club size and diminish participation rates in larger schools. Furthermore, national organizations like Key Club usually establish one chapter per school, inadvertently causing the chapter to become unwieldy in larger schools, further dissuading participation.

We propose extending the practice applied to courses to clubs. In large schools, it would be advantageous to establish multiple clubs of the same category, such as multiple chess clubs or Key Club chapters, each with independent club officers. This approach would provide leadership opportunities to a larger pool of students, who subsequently can work with the student community to boost participation rates. In addition, we recommend encouraging the creation of similar clubs with slight variations in focus. Notably, large schools successful in offering many clubs tend to feature numerous community service clubs, each centered around a slightly distinct theme. Conversely, certain schools enforce policies discouraging the establishment of new clubs similar to existing ones, inadvertently leading to the overmanning effect and reduced participation.

Our recommendation takes a more moderate stance compared to Leithwood and Jantzi's recommendation to cap secondary schools at 1,000 students or fewer (Leithwood & Jantzi, 2009). Fourteen years after their initial proposal, schools with 1,000 or fewer students host only 34% of the total student population (Figure 1), and the trend of school consolidation remains unaltered. We suggest working within the existing school structure and introducing smaller, immediately actionable changes with substantial potential, instead of waiting for major school demographics to shift, which could take decades to occur.

Besides our more actionable approach to coping with the hard-to-change school size factor, in general, we strongly advocate for individual schools taking initiatives to improve club offerings despite the constraints of school

demographics. Specifically, Figure 12 shows that the top 1/4 schools have 3.2 times more clubs than the bottom 1/4 schools after school demographics are taken into account, emphasizing the pivotal role of individual school initiatives. Besides these statistics, there are also concrete examples. One inspiring example is Syracuse Academy of Science Charter School (Syracuse, 2023). Despite having only 286 students, with 75% of them receiving free or reduced-price lunch, the school boasts an impressive 29 clubs, which is several times higher than that of similar schools. The school's focus on science is reflected in its numerous STEM-related clubs, but it also offers a diverse range of clubs encompassing humanity, charity, arts, and hobbies.

The statement from Principal Corey Tafoya of Woodstock High School, which had successfully improved the student participation rate in extracurricular activities by over 400% in five years, encapsulates our recommendation most effectively: "If we have six or seven students interested in something, we'll start a new club. We want students to find a reason to get up and come to school. Whatever trips their trigger is what our teachers and administration are willing to do" (Reeves, 2008).

## 5 Limitations and Future Work

One limitation of this research is that it collects samples from school websites, which could lead to biased outcomes, as it excludes schools that do not publish club data online. Despite this limitation, it is considered an acceptable tradeoff for the first attempt, given the constraints of alternative methods for gathering comprehensive club data. For instance, conducting in-person surveys would likely introduce bias toward schools willing to participate in this type of study, and achieving widespread coverage across hundreds of schools dispersed throughout the United States would be challenging. Other limitations include the exclusion of varsity sports clubs and private schools from this study, which are subjects for future research.

## 6 Related Works

The past research has presented compelling evidence of a strong association between positive youth development and active involvement in extracurricular activities, as supported by various studies (Mkude & Mubofu, 2022;

Leksuwankun, Benítez, Albertos, & Lara, 2023; Balaguer et al., 2020a; Busseri et al., 2006; Gilman et al., 2004; Eccles et al., 2003; Marsh & Kleitman, 2002; Darling et al., 2005; Eccles & Templeton, 2002; Mahoney et al., 2005; Zaff et al., 2003; Lerner et al., 2005; Gardner et al., 2008; Reeves, 2008; Fredricks & Eccles, 2006; Fredricks & Eccles, 2005; Peck et al., 2008). In addition, several surveys have summarized the effects of participation in extracurricular activities (Seow & Pan, 2014; Feldman & Matjasko, 2005; Holland & Andre, 1987; Farb & Matjasko, 2012; Rahayu & Dong, 2023). Specifically, it has been shown that positive parenting is associated with success in extracurricular activities and the development of personality traits (Balaguer et al., 2020b). Finally, a prominent taxonomy, known as “the five Cs,” attributes the positive outcomes to the beneficial effect of organized extracurricular activities on five key areas of youth development: competence, confidence, connection, character, and caring (Lerner et al., 2005).

It is important to note that, while there is evidence of a strong association between positive youth development and extracurricular activities, some studies argue that factors other than the quantity and variety of activities offered by a school – such as peer relations, sense of belonging, perceived support, and school climate – play a more important role (Balaguer, Benítez, de la Fuente, & Osorio, 2022; Hamlin, 2021; Berhanu & Sewagegn, 2024).

Concerns about the over-scheduled child problem have been raised in some studies (Rosenfeld & Wise, 2010). However, Mahoney et al. conducted an extensive survey and supported promoting participation in extracurricular activities, as they found limited empirical support for the overscheduling hypothesis and consistent evidence for the positive youth development perspective (Mahoney et al., 2006).

Some research suggests that the benefits of extracurricular activities may depend, in part, on the type of activities in which youth participate (Marsh & Kleitman, 2002; Larson, Hansen, & Moneta, 2006). Fredricks & Eccles (2006) examined the effect of the total number and breadth of participation in activities on youth development. These studies assume that a reasonable number and variety of extracurricular clubs are available to students. Related to this assumption, this research examines factors influencing the availability of high school clubs.

Barker and Gump’s study on school size and available extracurricular activities is relevant to this research (Barker & Gump, 1964). However, their work was limited to data from a small number of schools in a specific region (13 high schools in Eastern Kansas) and was conducted about half a century ago. In contrast, our modern research is more comprehensive, encompassing 229 schools spread

across the United States, and considering factors beyond just school size.

McNeal’s study (McNeal, 1999) is also related to this research, as it investigated the effect of school size and pupil-to-teacher ratio on student participation in high school extracurricular activities. However, in terms of the student–school relationship, this research focuses on the supply side of extracurricular clubs offered by high schools, in contrast to McNeal’s emphasis on the demand side, specifically student participation. Similarly, some work primarily focuses on factors affecting non-participants, such as lower socioeconomic status, lower grades, and larger schools (Feldman & Matjasko, 2007).

Various studies have discussed the effects of large school sizes, including their affect on student indiscipline (Haller, 1992), dropout rates (Alspaugh, 1998), voluntary participation (Schoggen & Schoggen, 1988), social participation (Morgan & Alwin, 1980), and social networks (Schaefer, Simpkins, Vest, & Price, 2011). Stevens and Peltier conducted a literature review and found support for the claim that students in smaller schools are more actively involved in extracurricular activities than students in larger schools (Stevens & Peltier, 1994). Our quantification of the  $\frac{1}{\sqrt{X_{\text{enroll}}}}$  pattern of clubs per student helps shed light on the root cause of many observations above.

## 7 Conclusion

By using data collected from hundreds of American high schools, we have studied the factors that impact schools’ club counts, which ultimately affect students’ participation rates in club activities. Our findings are summarized below.

First, the number of clubs offered by a school follows the trend of  $\sqrt{X_{\text{enroll}}}$ , where  $X_{\text{enroll}}$  is the school’s student enrollment. The correlation between school size and the number of clubs has a coefficient of 0.63 and  $p < 0.001$ . Consequently, the number of clubs per student follows the trend of  $\frac{\sqrt{X_{\text{enroll}}}}{X_{\text{enroll}}}$  (i.e.,  $\frac{1}{\sqrt{X_{\text{enroll}}}}$ ). This implies that although larger schools offer a wider array of club options, their clubs have either larger sizes or lower participation rates, or both, leading to reduced engagement. Although the previous research has extensively highlighted the negative effects of larger school sizes on student participation (Schoggen & Schoggen, 1988; Morgan & Alwin, 1980; Schaefer et al., 2011; Stevens & Peltier, 1994; Leithwood & Jantzi, 2009), to our knowledge, this study is the first to precisely quantify the  $\frac{1}{\sqrt{X_{\text{enroll}}}}$  pattern of clubs per student.

This quantification sheds light on the root cause of many observations documented in prior research.

Second, using the fraction of students in a school receiving free or reduced-price lunch, denoted as  $X_{\text{lunch}}$ , as an indicator of household income, schools in the higher-income group on average have 16.7 more clubs than those in the lower-income group. Furthermore, when  $X_{\text{lunch}}$  is less than 40%, an increase in  $X_{\text{lunch}}$  leads to a decrease in the number of clubs. However, when  $X_{\text{lunch}}$  exceeds 40%, a further increase in  $X_{\text{lunch}}$  no longer affects the number of clubs. For context, the average value of  $X_{\text{lunch}}$  for schools is 50%. This suggests that the advantages of having more clubs are primarily observed in higher-income schools, while average-income schools do not significantly differ from lower-income schools. More precisely, the nonlinear impact of  $X_{\text{lunch}}$  on the number of clubs can be effectively modeled using a logistic function.

Third, at first glance, there seems to be a positive correlation between the pupil-to-teacher ratio and the number of clubs, implying that a decrease in the number of teachers leads to an increase in the number of clubs. However, this counterintuitive outcome is primarily due to schools with higher pupil-to-teacher ratios mostly being large schools, which generally offer more clubs. Once the influence of school size is excluded, the number of clubs exhibits a negative correlation with the pupil-to-teacher ratio. Specifically, a controlled experiment shows that among schools of comparable sizes, schools with lower pupil-to-teacher ratios, on average, have 9.3 more clubs than those with higher pupil-to-teacher ratios.

Fourth, at first glance, racial demographics seem to affect the number of clubs, as they are correlated with school size and household income, which directly affect the number of clubs. However, after excluding their indirect influence through school size and household income, racial demographics by themselves do not significantly affect the number of clubs. Specifically, both  $t$ -test and multiple regression results independently confirm that, among schools with comparable student enrollments and household incomes, but significant differences in their racial demographics, no significant difference exists in their club offerings.

Fifth, our analysis reveals that for schools with identical demographic characteristics – such as school size, household income, and pupil-to-teacher ratio – the top 25% of schools offer 3.2 times more clubs than the bottom 25%. This notable difference underscores the crucial role of school initiatives in increasing club offerings.

In summary, the factors influencing the number of clubs offered by schools, ranked by impact, are school size, household income, and pupil-to-teacher ratio, while

racial demographics do not significantly affect club count. Finally, despite demographic constraints such as school size, schools still have ample opportunities to take the initiative to improve club offerings.

**Acknowledgments:** The authors are grateful for the reviewers' valuable comments that improved the manuscript.

**Funding information:** The authors state no funding involved.

**Author contributions:** All authors take responsibility for the entire content of this manuscript, have consented to its submission to the journal, reviewed all the results, and approved the final version of the manuscript. A.T. conceived the original idea for this study. A.T. and C.T. collected the data and performed the analysis. Y.W. proposed the multiple-regression method. All authors contributed to the writing and finalization of the manuscript.

**Conflict of interest:** The authors state no conflict of interest.

**Data availability statement:** The data will be provided upon request to the authors.

## References

- Alspaugh, J. W. (1998). The relationship of school-to-school transitions and school size to high school dropout rates. *The High School Journal*, 81(3), 154–160.
- Anjum, S. (2021). Impact of extracurricular activities on academic performance of students at secondary level. *International Journal of Applied Guidance and Counseling*, 2(2), 7–14.
- Balaguer, A., Benítez, E., Albertos, A., & Lara, S. (2020a). Not everything helps the same for everyone: Relevance of extracurricular activities for academic achievement. *Humanities and Social Sciences Communications*, 7(1), 1–8.
- Balaguer, A., Benítez, E., de la Fuente, J., & Osorio, A. (2022). Structural empirical model of personal positive youth development, parenting, and school climate. *Psychology in the Schools*, 59(3), 451–470.
- Balaguer, A., Orejudo, S., Rodríguez-Ledo, C., & Cardoso-Moreno, M. J. (2020b). Extracurricular activities, positive parenting and personal positive youth development. Differential relations amongst age and academic pathways. *Electronic Journal of Research in Educational Psychology*, 2(18), 179–206.
- Bandura, A. & Walters, R. H. (1977). *Social learning theory*. Englewood Cliffs, NJ: Prentice Hall.
- Barker, R. G., & Gump, P. V. (1964). *Big school, small school: High school size and student behavior*. Redwood City, CA: Stanford University Press.
- Berhanu, K. Z., & Sewagegn, A. A. (2024). The role of perceived campus climate in students' academic achievements as mediated by

- students engagement in higher education institutions. *Cogent Education*, 11(1), 2377839.
- Borrego, I., & Cuadrado, I. (2025). Extracurricular activities for adolescents to prevent social networks addiction. *American Research Journal of Humanities Social Science*, 8(2), 57–63.
- Buckley, P., & Lee, P. (2021). The impact of extra-curricular activity on the student experience. *Active Learning in Higher Education*, 22(1), 37–48.
- Busseri, M. A., Rose-Krasnor, L., Willoughby, T., & Chalmers, H. (2006). A longitudinal examination of breadth and intensity of youth activity involvement and successful development. *Developmental Psychology*, 42(6), 1313.
- Darling, N., Caldwell, L. L., & Smith, R. (2005). Participation in school-based extracurricular activities and adolescent adjustment. *Journal of Leisure Research*, 37(1), 51–76.
- Deci, E. L., & Ryan, R. M. (2012). Self-determination theory. *Handbook of Theories of Social Psychology*, 1(20), 416–436.
- Dewey, J., & Wheeler, J. E. (1913). *Interest and effort in education*. Boston: Houghton Mifflin.
- Eccles, J. S., Barber, B. L., Stone, M., & Hunt, J. (2003). Extracurricular activities and adolescent development. *Journal of Social Issues*, 59(4), 865–889.
- Eccles, J. S., & Templeton, J. (2002). Chapter 4: Extracurricular and other after-school activities for youth. *Review of Research in Education*, 26(1), 113–180.
- Farb, A. F., & Matjasko, J. L. (2012). Recent advances in research on school-based extracurricular activities and adolescent development. *Developmental Review*, 32, 1–48.
- Feldman, A. F., & Matjasko, J. L. (2005). The role of school-based extracurricular activities in adolescent development: A comprehensive review and future directions. *Review of Educational Research*, 75(2), 159–210.
- Feldman, A. F., & Matjasko, J. L. (2007). Profiles and portfolios of adolescent school-based extracurricular activity participation. *Journal of Adolescence*, 30(2), 313–332.
- Feraco, T., Resnati, D., Fregonese, D., Spoto, A., & Meneghetti, C. (2023). An integrated model of school students' academic achievement and life satisfaction. linking soft skills, extracurricular activities, self-regulated learning, motivation, and emotions. *European Journal of Psychology of Education*, 38(1), 109–130.
- Fredricks, J. A., & Eccles, J. S. (2005). Developmental benefits of extracurricular involvement: Do peer characteristics mediate the link between activities and youth outcomes? *Journal of Youth and Adolescence*, 34, 507–520.
- Fredricks, J. A., & Eccles, J. S. (2006). Extracurricular involvement and adolescent adjustment: Impact of duration, number of activities, and breadth of participation. *Applied Developmental Science*, 10(3), 132–146.
- Gardner, M., Roth, J., & Brooks-Gunn, J. (2008). Adolescents' participation in organized activities and developmental success 2 and 8 years after high school: Do sponsorship, duration, and intensity matter? *Developmental Psychology*, 44(3), 814.
- Gilman, R., Meyers, J., & Perez, L. (2004). Structured extracurricular activities among adolescents: Findings and implications for school psychologists. *Psychology in the Schools*, 41(1), 31–41.
- Haller, E. J. (1992). High school size and student indiscipline: Another aspect of the school consolidation issue? *Educational Evaluation and Policy Analysis*, 14(2), 145–156.
- Hamlin, D. (2021). Can a positive school climate promote student attendance? evidence from new york city. *American Educational Research Journal*, 58(2), 315–342.
- Heath, R. D., Anderson, C., Turner, A. C., & Payne, C. M. (2022). Extracurricular activities and disadvantaged youth: A complicated but promising story. *Urban Education*, 57(8), 1415–1449.
- Holland, A., & Andre, T. (1987). Participation in extracurricular activities in secondary school: What is known, what needs to be known? *Review of Educational Research*, 57, 437–466.
- Kleinert, E. J. (1969). Effects of high school size on student activity participation. *The Bulletin of the National Association of Secondary School Principals*, 53(335), 34–46.
- Larson, R. W., Hansen, D. M., & Moneta, G. B. (2006). Differing profiles of developmental experiences across types of organized youth activities. *Developmental Psychology*, 42(5), 849–863.
- Lave, J. (1991). *Situated learning: Legitimate peripheral participation*. Cambridge: Cambridge University Press.
- Leithwood, K., & Jantzi, D. (2009). A review of empirical evidence about school size effects: A policy perspective. *Review of Educational Research*, 79(1), 464–490.
- Leksuwankun, S., Dangprapai, Y., & Wangsaturaka, D. (2023). Student engagement in organising extracurricular activities: Does it matter to academic achievement? *Medical Teacher*, 45(3), 272–278.
- Lerner, R. M., Lerner, J. V., Almerigi, J. B., Theokas, C., Phelps, E., Gestsdottir, S., ..., von Eye, A. (2005). Positive youth development, participation in community youth development programs, and community contributions of fifth-grade adolescents: Findings from the first wave of the 4-h study of positive youth development. *The Journal of Early Adolescence*, 25(1), 17–71.
- MacKinnon, J. G., & White, H. (1985). Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties. *Journal of Econometrics*, 29(3), 305–325.
- Mahoney, J. L., Harris, A. L., & Eccles, J. S. (2006). Organized activity participation, positive youth development, and the over-scheduling hypothesis. *Social Policy Report*, 20(4), 3–30.
- Mahoney, J. L., Larson, R. W., & Eccles, J. S. (2005). *Organized activities as contexts of development: Extracurricular activities, after school and community programs*. London, England: Psychology Press.
- Marsh, H., & Kleitman, S. (2002). Extracurricular school activities: The good, the bad, and the nonlinear. *Harvard Educational Review*, 72(4), 464–515.
- McNeal Jr, R. B. (1999). Participation in high school extracurricular activities: Investigating school effects. *Social Science Quarterly*, 80(2), 291–309.
- Mkude, M., & Mubofu, C. (2022). Extracurricular activities in the broader personal development: Reflections from youth in public secondary schools. *Research Ambition: An International Multidisciplinary e-Journal*, 6(IV), 1–5.
- Morgan, D. L., & Alwin, D. F. (1980). When less is more: School size and student social participation. *Social Psychology Quarterly*, 43(2), 241–252.
- Mukesh, H. V., Acharya, V., & Pillai, R. (2023). Are extracurricular activities stress busters to enhance students well-being and academic performance? Evidence from a natural experiment. *Journal of Applied Research in Higher Education*, 15(1), 152–168.
- NCES. (2023). National Center for Education Statistics. <https://nces.ed.gov/ccd/elsi/tableGenerator.aspx>.
- Peck, S. C., Roeser, R. W., Zarrett, N., & Eccles, J. S. (2008). Exploring the roles of extracurricular activity quantity and quality in the educational resilience of vulnerable adolescents: Variable-and pattern-centered approaches. *Journal of Social Issues*, 64(1), 135–156.
- Rahayu, A. P., & Dong, Y. (2023). The relationship of extracurricular activities with students' character education and influencing

- factors: a systematic literature review. *Al-Ishlah: Jurnal Pendidikan*, 15(1), 459–474.
- Reeves, D. B. (2008). The learning leader/the extracurricular advantage. *Learning*, 66(1), 86–87.
- Rosenfeld, A., & Wise, N. (2010). *The over-scheduled child: Avoiding the hyper-parenting trap*. New York: St. Martin's Griffin.
- Schaefer, D. R., Simpkins, S. D., Vest, A. E., & Price, C. D. (2011). The contribution of extracurricular activities to adolescent friendships: New insights through social network analysis. *Developmental Psychology*, 47(4), 1141–1152.
- Schoggen, P., & Schoggen, M. (1988). Student voluntary participation and high school size. *The Journal of Educational Research*, 81(5), 288–293.
- Seow, P.-S., & Pan, G. (2014). A literature review of the impact of extracurricular activities participation on students' academic performance. *Journal of Education for Business*, 89(7), 361–366.
- Statsmodels. (2023). <https://www.statsmodels.org/>.
- Stevens, N. G., & Peltier, G. L. (1994). A review of research on small-school student participation in extracurricular activities. *Journal of Research in Rural Education*, 10(2), 116–120.
- Syracuse. (2023). Syracuse Academy of Science Charter School, NCES school ID 360010405549, list of clubs: <https://www.sascs.org/students-parents/student-activities>.
- Wigfield, A., & Eccles, J. S. (2000). Expectancy-value theory of achievement motivation. *Contemporary Educational Psychology*, 25(1), 68–81.
- Zaff, J. F., Moore, K. A., Papillo, A. R., & Williams, S. (2003). Implications of extracurricular activity participation during adolescence on positive outcomes. *Journal of Adolescent Research*, 18(6), 599–630.