

# Case Study 3

## Elementary Education

Sam Lee, Ty Hawkes

This study examines the determinants of standardized test performance across California school districts, using a comprehensive dataset that includes key socioeconomic variables such as district income, percentage of English learners, and computer availability. We employ two modeling approaches: a multiple linear regression with a quadratic term for income to capture non-linear effects, and a generalized additive model (GAM) to allow for flexible functional forms. While both models demonstrate strong predictive performance, the linear regression is favored for its interpretability and parsimony. Our analysis highlights a significant negative impact of English learner concentration on test scores and provides robust evidence of diminishing returns to income on academic performance.

### Introduction

Research shows that strong academic performance during a child's elementary school years is a strong predictor of their successes later in life. Understanding which things are related to a student's academic performance in elementary school can help educators, administrators, and government leaders make informed decisions that positively affect rising generations. In this analysis, we hope to inform school officials and policy makers about elementary test scores and the factors that may affect them. We will study the state-wide standardized test scores of various school districts in California and examine several factors associated with an increase or decrease in overall test scores.

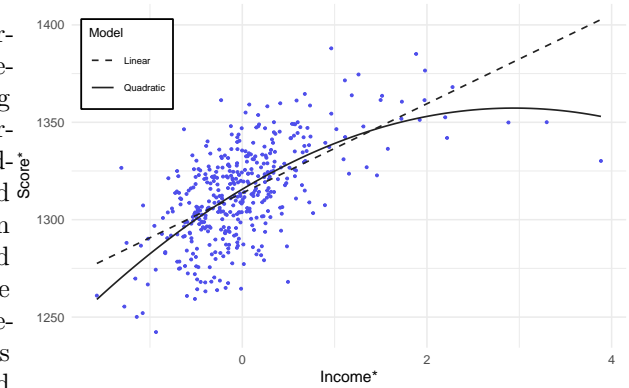


Figure 1: Non-linear Relationship between Income and Score

After conducting an exploratory data analysis, we found that there are two potential issues that could affect our primary analysis. Some of the factors in this study, like the district average income and the percentage of students who qualify for reduced-price lunch, are closely related to each other. If ignored, it can be difficult to determine the relationship between these factors and test scores. To avoid this, we will run tests to evaluate the severity of the issue, and remove factors if they pose a big enough problem. In addition to this, district average income appears to have a non-linear<sup>1</sup> (See Figure 1) relationship with test scores. This reduces the accuracy of our results if not accounted for. To solve this issue, we will add another factor to

<sup>1</sup>We visually model the relationship between *Income* and *Score* by partialling out the regressor (*Income*) and the regressand (*Score*). Thus, we compute  $Income^*$  as  $Income - Z((Z'Z)^{-1}Z'Income)$ , where  $Z$  are the set of covariates excluding *Income*; the set of covariates are an  $n \times (k - p)$  matrix, where  $p$  is the number of covariates that are “partialled out” (the *Income* in this case). Similarly, we compute  $Score$  as  $Score - Z((Z'Z)^{-1}Z'Score)$ . Hence, Figure 1 represents the non-linear effect that income has on *Score*, holding all else constant. Note that in  $Income^*$  represents the scaled *Income*.

our analysis that will improve our model accuracy and provide additional insight, at the expense of some interpretability.

## Methodology

We first fit a LASSO regression model with a second-degree polynomial term included for *Income* to assess variable selection. Through this process, we eliminated using *STratio* (student-to-teacher ratio) as a predictor as the LASSO model shrunk *STratio* to zero. We used the LASSO-selected covariates (all of which have descriptions summarized in Table 1) for all models going forward.

We first propose a multiple linear regression model with an added second degree polynomial term for income. This model is a good candidate because it accounts for the non-linearity present in *Income*. More importantly, it will allow us to evaluate the relationships between test scores and various factors due to its parsimony. Despite these strengths, this model may not be as predictive of student scores as other models. Additionally, this model will only fit well if the relationship between test scores and income is quadratic and the other relationships between the factors and *Score* holds linearly.

Next, we propose a Generalized Additive Model (henceforth known as GAM). This model is a good candidate because it will also account for the non-linearity of the data with smoothing techniques. One advantage to using this model over the linear regression model is that GAM regression can model complex, non-linear relationships that may not be captured well with polynomial expansions, possibly leading to a better fit. Due to this flexibility, however, GAM regression does not provide estimates for the effect size that each factor has on test scores. Still, we are able to accomplish the goals of this analysis with this type of model because GAM regression allows us to determine the statistical significance of these relationships, and visually interpret their direction.

Both models assume independent and Normally distributed data. The linear model imposes stricter assumptions, including homoskedastic variance and a linear relationship between the covariates and response. In contrast, GAM, though non-parametric and flexible in capturing non-linear effects, also assumes homoskedastic errors under the Normality assumption.

## Model Evaluation

We first tuned our GAM model by selecting optimal hyperparameters for each covariate in the model. This process was achieved through a *randomized grid search* over the parameter space of interest. We chose an optimal basis function and the optimal number of “knots” for each factor and cross-validated each selection of hyperparameters through k-fold (with  $k = 5$ ) validation. The final form of our GAM model<sup>2</sup> can be represented by Equation 1:

$$y_i = \beta_0 + \sum_{j=1}^{13} \alpha_j \phi_j(\text{Lunch}_{i1}) + \sum_{j=1}^4 \beta_j \psi_j(\text{Computer}_{i2}) + \sum_{j=1}^{14} \gamma_j \psi_j(\text{Expenditures}_{i3}) + \sum_{j=1}^{19} \delta_j \phi_j(\text{Income}_{i4}) + \sum_{j=1}^{20} \zeta_j \phi_j(\text{English}_{i5}) + \varepsilon_i \quad (1)$$

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

We evaluated both our GAM regression model and Linear regression model on their in-sample and out-of-sample performance measures. The out-of-sample RMSE for each model was evaluated using k-fold cross-validation (using  $k = 20$ ). These results are summarized in Table 3. Both models performed exceptionally well when cross validated.

---

<sup>2</sup>Nomenclature:  $\alpha_j, \beta_j, \gamma_j, \delta_j$ , and  $\zeta_j$  represent spline coefficients.  $\phi_j(\cdot)$  represent thin-plate spline basis functions, and  $\psi_j(\cdot)$  represent cubic basis functions.

Our linear regression model had an adjusted R squared value of 0.7877 and our GAM regression model had an adjusted R squared value of 0.7796. Both models fit the data well, with the linear regression achieving a higher adjusted R-squared due to its parsimonious fit. Because both models showed similar predictability, we ultimately chose to use our linear model in favor of its superior interpretability. Our linear model can be represented by the solution to the linear combination shown in Equation (2) below.

$$y_i = \beta_0 + \beta_1 \text{Lunch}_i + \beta_2 \text{Computer}_i + \beta_3 \text{Expenditure}_i + \beta_4 \text{English}_i + \beta_5 \text{Income}_i + \beta_6 \text{Income}_i^2 + \epsilon_i \quad (2)$$

$$\epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

To ensure the validity of our linear regression model, we evaluated its adherence to the fundamental assumptions of linear regression. Figure 2 assesses the linearity between the covariates and response with added variable plots. The linearity in each plots provides strong evidence that this assumption is met. The assumption of independence is assumed since the observations were collected without any known time-based or spatial dependencies. The scale-location plot in Figure 3 assesses the homoscedasticity of our residuals. The relatively flat line indicates a constant variance in our errors. Lastly, the Q-Q plot in Figure 3 assesses the assumption that our errors are normally distributed. The flat slope indicates that the errors of our data follow the expected quantile measurements in a normal distribution.

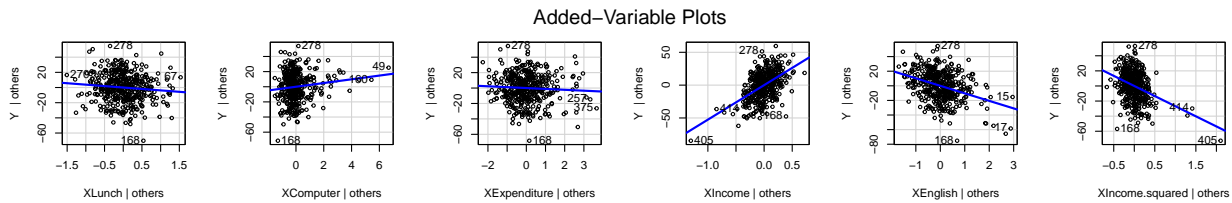


Figure 2: Added variable plots to assess linearity of linear model

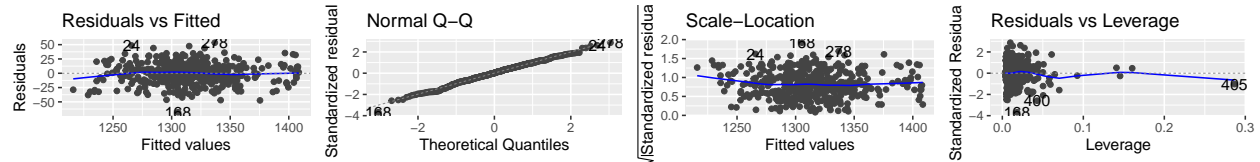


Figure 3: Model diagnostics to assess linear model assumptions of Equation (1)

## Results

We present our results in terms of the estimation of Equation 2 by showing the estimates of coefficients on the scaled factors below in Table 1—in other words, each estimated coefficient represents a *relative* effect on *Score* as it relates to the other coefficients. Hence, we compare the magnitude in the coefficients (and their respective standard errors) to assess which factors contribute most significantly.

Table 1: Regression Results

Variable	Description	Estimate	95% Confidence Interval
Intercept		1313.666***	(1311.872, 1315.46)
Computer	Number of Computers	2.458*	(0.48, 4.436)
English	Percent of English learners	-10.317***	(-12.906, -7.728)
Expenditure	Expenditure per student	-1.145	(-3.121, 0.831)
Income	District average income (in USD 1,000)	52.568***	(43.849, 61.287)
Income <sup>2</sup>	Income squared	-26.815***	(-34.285, -19.346)
Lunch	Percent qualifying for reduced-price lunch	-3.777*	(-7.461, -0.094)

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1 | Multiple R-squared: 0.7907, Adjusted R-squared: 0.7877  
Residual standard error: 18.7 on 413 degrees of freedom | F-statistic: 260.1 on 6 and 413 DF, p-value:  $< 2.2 \times 10^{-16}$

We return to Figure 1, and address the non-linearities in *Score* with respect to *Income*. We assess whether the quadratic term, *Income*<sup>2</sup> contributes significantly to explaining the variability in *Score* by running an analysis of variance test (ANOVA) on two models: The first being the empirical model established by Equation 2; The second, like unto the first, omitting the quadratic term. The results of this test are summarized in Table 2. We reject the null hypothesis that there is no difference between the two models. In other words, the quadratic term cannot be omitted without sacrificing predictability in *Score*. Since the coefficient on *Income*<sup>2</sup> is significantly negative, this provides significant evidence to suggest that as income increases, the total effect that income has on increasing *Score*, decreases. This supports the *diminishing marginal returns* hypothesis.

Table 2: Analysis of Variance Table

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	413	144446				
2	414	161863	-1	-17417	49.799	$7.238 \times 10^{-12}$

Model 1:  $Y = X\beta + \epsilon$   
Model 2:  $Y = X_{-Income^2}\beta_{-Income^2} + \epsilon$

To assess predictability, we compare our linear model to GAM (as proposed earlier) and local linear regression with a quadratic term (LOESS). Similar to our tuned GAM model, we tuned our LOESS model through k-fold cross validation<sup>3</sup>. We performed k-fold cross validation using  $k = 20$ , and summarize the following results below as the out-of-sample predictive fit for each model:

Table 3: Out of Sample RMSE for each Predictive Model

Linear Model	LOESS	GAM
18.51708	18.23348	18.53671

We also summarize this visually in Figure 4 by modeling how well each model predicts using *Income* as a predictor. We visually compare our quadratic model (Equation 2) to LOESS by fitting both on the partialled-out *Income* (since the LOESS was fitted on the dimensionally-reduced covariate matrix) and *Score*. Similarly, we compare our quadratic model to a more flexible GAM model by fitting both to the original scaled *Income* while holding all other covariates constant. We maintain that the quadratic model proves as the best predictive and parsimonious model for this set of data.

<sup>3</sup>Our LOESS model was tuned through a randomized grid search algorithm, where we optimized over a space of *span* (the percentile of data used for each ‘local’ regression) parameters in the set (0.1, 1). Through five-fold cross-validation, we found the optimal span to be 0.9919. We note that since this was optimized over a random selection over the parameter space from (0.1, 1), the theoretical optimal may approach 1. To optimize algorithmic performance, we performed dimension reduction by partialling out the set of regressors ( $Z$ ) on both (scaled) *Income* and *Score*.

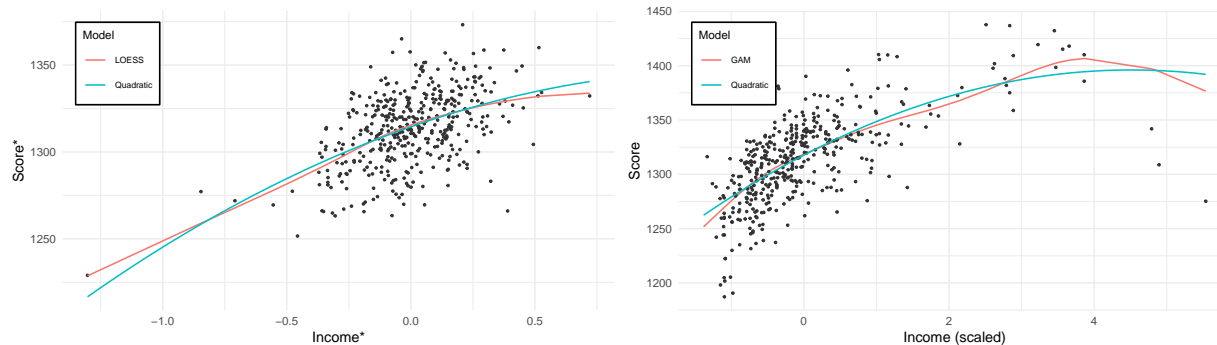


Figure 4: Model Prediction on *Score* with respect to *Income*

We now turn to address the more precarious causal questions we have been presented with. Table 1 suggests that the factor, *English*, is significantly negative—in other words, for every one percent increase in the number of English learners within a school district, our model suggests that the average score for that district will decrease by -0.5642 (unscaled), on average. Hence, learning English as a second language may be a barrier to learning if<sup>4</sup> we do not suspect there are other unobserved covariates influencing the percent of English learners that are *also* correlated with *Score*. Unobserved factors such as school funding, parental education levels, teacher quality, and immigration status are all factors that confound the causal inference assumption. Thus, we caution from making any binding causal claims on this issue.

The number of computers per district is the strongest positive predictor of *Score*. However, to make causal claims, we recommend a randomized experiment: assign varying numbers of computers to different schools, ensuring the number of computers is not correlated with other factors like teacher quality or socioeconomic status. An intervention group would receive additional computers, while a control group would maintain the current number, isolating the effect of *Computers* on *Score*.

## Conclusion

The goal of this study was to evaluate the relationships between the test scores of elementary students in California and their respective school circumstances, to shed some light on possible administrative methods to improve student learning. We fit a linear regression model to our data and found that there is evidence of diminishing returns on extra curricular activities on student learning. We also found that schools with a higher percentage of students learning English as a second language tend to have a lower average score on the Stanford 9 test. Lastly, we found that among things that a school district can control, a higher number of computers present at a school was correlated with higher test scores.

The biggest shortcoming of this study is it's observational nature. There is no evidence that any of our covariates were assigned randomly, thus we cannot make any causal statements to confidently state whether any of these factors actually have an effect on test scores. The effect of computers on student learning should be further investigated, possibly using randomized controlled experiments to allow for a causal claim that computers actually do boost student learning. Additionally, more studies and/or experiments should be done to better understand how English learners are possibly at a disadvantage academically, and to better understand how we can help level the playing field.

## Teamwork

<sup>4</sup>For a sufficient sample size  $n$ , if the endogeneity assumption holds, that is, if we have sufficient evidence to believe that  $\mathbb{E}[\varepsilon|X] = 0$ , then a causal claim may be warranted. However, we caution against this due to confounding factors in the data as we mentioned.

Sam spearheaded the hyper parameter tuning, designed the plots, wrote the results section of the paper, and beautified many aspects of the final report. Ty spearheaded the EDA and modeling, verified assumptions, and wrote the introduction, methodology, and conclusion sections of the final report.