
Data Driven Learning Intervention: Novel Approach and Case Study

Educational Researcher
XX(X):1–23
©The Author(s) 2024
Reprints and permission:
sagepub.co.uk/journalsPermissions.nav
DOI: 10.1177/ToBeAssigned
www.sagepub.com/

SAGE

Abstract

In this paper, we present a methodology and path forward for data-driven quantitative analysis to support qualitative findings when analyzing large open-source datasets. We focus on the task of learning interventions and policies and use nine publicly accessible datasets to understand and mitigate factors contributing to learning loss and the practical learning recovery measures in Texas public school districts after the recent school closures. The data came from the Census Bureau 2010, USAFACTS, Texas Department of State Health Services (DSHS), the National Center for Education Statistics (CCD), the U.S. Bureau of Labor Statistics (LAUS), and three sources from the Texas Education Agency (STAAR, TEA, ADA, ESSER). We demonstrate a novel data-driven approach to discover insights from an extensive collection of heterogeneous public data sources. For the pandemic school closure period, the mode of instruction and prior score emerged as the primary resilience factors in the learning recovery intervention method. Grade level and census community income level are the most influential factors in predicting learning loss for both math and reading. We demonstrate that data-driven unbiased data analysis at a larger scale can offer policymakers an actionable understanding of how to identify learning-loss tendencies and prevent them in public schools.

Introduction

Learning loss, within the context of education, can be defined as the depletion or regression of previously attained or expected knowledge and competencies. COVID-19 also had an impact on teacher preparation Choate et al. (2021). As an example, the recent COVID-19 pandemic forced many schools to close, and the global consequences of a five-month closure of schools were equated to less than \$10 trillion monetary loss as 43 million students were affected by the school closures OECD (2021). The school closures have led to learning loss among students and higher teacher attrition

Courtemanche et al. (2021); Zamarro et al. (2022). The learning loss percentage in some countries was estimated from 0.08 to 0.29 based on the public data Maldonado and De Witte (2022). In the United States, the school closure and subsequent reopening were uneven as there was no consensus. Thus, the learning loss was not uniform across the states, as documented for Virginia, Maryland, Ohio, and Connecticut in Halloran et al. (2021). Individual studies of Rhode Island and North Carolina education data provided estimations of the learning losses and recovery of Public Instruction (2022). In summary, more definitive data-driven factors contributing to the learning recovery in the studies mentioned above remain elusive and diverse Betebenner et al. (2021). This complexity posed challenges for policymakers in the other states in determining the learning intervention measures based on their data.

In this paper, we examine if the intervention model matches the publicly available data in Texas. The Texas Education Agency published a report documenting the 4% **Loss** in reading and 15% **Loss** in math on the STAAR exam and how the negative impact of COVID-19 erased years of improvement in reading and math Agency (2021); , 2022b (TEA). If the apparent factors (census data), location (urban vs. rural), and standardized exam data are good predictors of learning loss. Next, we offer to pinpoint resilient factors contributing to learning recovery within Texas public schools.

Our approach is novel in that it integrates data science methodologies with educational policy analysis, offering a data-driven perspective to inform decision-making processes. By aggregating and analyzing data from nine sources, we aim to identify what specific factors were most important for the schools to experience significant learning loss and learning recovery. We propose improved automated attribute importance analysis to understand various parameters, including consensus information, demographics of public school districts, instructional modalities, socioeconomic indicators such as income levels, urban or rural settings, student attendance rates, county infection rates, and unemployment statistics, among numerous other factors spanning the years 2019, 2021, and 2022.

We uncover fascinating data-driven indicators. The most resilient factor influencing learning loss in the district is how early or late the students go back to in-person learning. The size and location of a district, along with the amount of money in the area and the Elementary and Secondary School Emergency Relief Fund received, play critical roles in the recovery process. The results identify the significance of various factors in promoting learning recovery in math and reading, highlighting the importance of considering a district's economic status, size, locale, demographics, and funding. The remainder of this paper is structured as follows: Section reviews pertinent literature, Section outlines the research design, Section describes data gathering and preparation; Results presents our findings; and Section discusses the implications and suggests directions for future research.

Related Work

In this section, we will focus on (1) quantitative research and machine learning tools to gain insight from the data on the relationship with the outcome without overfitting the features to the data or (2) the directions for selecting machine learning models for predicting learning loss with tabular data. The most popular machine

learning (ML) techniques (logistic regression, support vector machines, Bayesian belief network, decision trees, and neural network) for data in the wild generally offer an excellent classification accuracy above 70% for simple classification tasks Cardona et al. (2020). From a data science perspective, it's critical to refine the selection of modeling approaches. Excessive reliance on feature engineering may result in less-than-optimal outcomes when translating domain-specific data Baashar et al. (2021). Machine Learning methods such as deep neural networks (DNN), decision trees, Support Vector Machines (SVM), and K-nearest neighbor (KNN) are favored methods for predicting student academic performance Rao et al. (2019). Also, the demographic, educational, familial/personal, and internal assessment factors emerged as the standard resource for data-driven student performance evaluation across various metrics such as classroom performance, grade levels, and standardized tests performance Baashar et al. (2021). A large-scale data science study examined the Big Fish Little Pond Effect (BFLPE), which describes how individuals often feel better about themselves when they excel in a less competitive environment rather than being average in a highly competitive one, across 56 countries for fourth-grade math and 46 countries for eighth-grade math. This analysis drew on extensive data from the Trends in International Mathematics and Science Study (TIMSS) and employed a straightforward statistical approach Wang (2020).

Recent research indicates that state-of-the-art machine learning techniques for tabular data surpass existing methods and exhibit less sensitivity to input bias and noise compared to Deep Neural Networks (DNN) Yan (2021). These findings suggest that for tasks involving tabular data, machine learning models such as gradient boosting machines and random forests may be more effective and robust choices than DNNs.

State-of-the-art gradient-boosted decision trees (GBDT) models such as XGBoost Chen and Guestrin (2016), LightGBM Ke et al. (2017), and CatBoost Prokhorenkova et al. (2018) are the most popular models of choice when it comes to tabular data due to their superior performance, scalability, and ability to handle various types of data and feature interactions effectively. These models are highly efficient in capturing complex patterns and relationships within tabular datasets, often leading to better predictive accuracy compared to other machine learning algorithms. In recent years, deep learning models have emerged as state-of-the-art techniques on heterogeneous tabular data: TabNet Arik and Pfister (2021), DNF-Net Abutbul et al. (2020), Neural Oblivious Decision Ensembles (NODE) Popov et al. (2019), and TabNN Ke et al. (2019). Although papers have proposed that these deep learning algorithms outperform the GBDT models, there is no consensus that deep learning exceeds GBDT on tabular data because standard benchmarks have been absent. There is also a shortage of open-source implementations, libraries, and their corresponding APIs for deep learning Joseph (2021); Schwartz-Ziv and Armon (2022). Recent studies provide competitive benchmarks comparing GBDT and deep learning models on multiple tabular data sets Schwartz-Ziv and Armon (2022); however, all of these benchmarks indicate that there is no dominant winner, and GBDT models still outperform deep learning in general. The studies suggest developing tabular-specific deep learning models such that tabular data modalities, spatial and irregular data due to high-cardinality categorical features, missing values, and uninformative features cannot guarantee the same prediction power

as deep learning obtains from homogeneous data, including images, audio, or text Grinsztajn et al. (2022).

Methodology

The methodology employed in this study aims to systematically uncover factors contributing to both learning loss and learning gain among students. By leveraging advanced educational data science techniques, the study utilizes comprehensive datasets from Texas public schools to analyze trends and patterns in math and reading scores. This methodology integrates statistical modeling, machine learning algorithms, and data visualization tools to identify critical variables influencing student academic performance over time. By examining multiple academic years, including the 2021-2022 data, the study ensures robustness and reliability in its findings, thereby offering valuable insights into educational outcomes and informing targeted interventions.

Attribute Importance Scoring

First, we introduce an innovative approach to identifying critical features from a vast array of potential factors. Heterogeneous data tends to have overlapping information mixed with numerical and categorical data. Logistic Regression coefficients for the actual data often randomly select one out of multiple correlated columns and are not robust enough for the noisy multi-source data analysis Kim et al. (2013). We propose to contrast filter, embedded, and wrapper methods for feature importance and propose a novel aggregation technique for robustness.

Several distinct (ten with variations) algorithms for automated feature selection are evaluated to assess the effectiveness of these techniques, along with a collection of interpretative methods for analyzing feature importance. These measures aim to mitigate issues associated with "Garbage In, Garbage Out (GIGO)" and trivial modeling. We construct a quasi-orthonormal attribute space by distilling and aggregating highly correlated features.

We wanted to avoid artificial weighting of the features in the modeling step, so we utilized this correlation filtering in this section to aggregate linearly related features in our data set into one attribute. To this end, we have expanded several categorical features to multiple binary features as we found that numerous separate categories capture highly overlapping data. The Pearson correlation coefficient ρ measures the linear relationship between two normally distributed variables and is defined in Equation 1:

$$\rho = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad (1)$$

The $\text{cov}(X, Y)$ represents the covariance between variables X and Y , while σ_X and σ_Y are the standard deviations of X and Y respectively. Pearson's correlation coefficient estimate r , also known as a "correlation coefficient," for attribute feature vectors $x = (x_1, \dots, x_n)$ with mean \bar{x} and $y = (y_1, \dots, y_n)$ with mean \bar{y} , is obtained via a Least-Squares fit, as defined in Equation 2.

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (2)$$

The \bar{x} and \bar{y} represent the means of vectors x and y respectively. A value of 1 represents a perfect positive relationship, -1 is a perfect negative relationship, and 0 indicates the absence of a relationship between variables. We use features with high correlation coefficients to aggregate them into one attribute, as they are linearly dependent on each other. Eventually, we could keep one attribute, the most highly correlated to our label, of those overlapping features in our analysis. Then, we can combine all binary dummy-coded variables from related categories as a set in variable selection. This approach thus reduces an attribute dimension that provides better interpretability of our attribute set and its importance. Here, we modify ten distinct approaches from filter methods, embedded methods, and wrapper methods sets to identify and assess the features influencing our prediction models. Each technique aims to select feature sets with minimal redundancy and maximal relevance, resulting in either a chosen set of features or a score indicating feature importance.

Table 1. Resilience Factors and Methods with Abbreviations

Abbreviation	Category	Full Term/Method
LI	Resilience Factor	Low Income
ATT	Resilience Factor	Attendance
DEM	Resilience Factor	Demographics
R/E	Resilience Factor	Race/Ethnicity
CC	Resilience Factor	County COVID
DM	Resilience Factor	District Makeup
MOI	Resilience Factor	Mode of Instruction
TST	Resilience Factor	Testing
PS	Resilience Factor	Prior Score
LOC	Resilience Factor	Locale
RFE RF	Feature Method	Recursive Feature Elimination - Random Forest
RFE RR	Feature Method	Recursive Feature Elimination - Ridge Regression
VT	Feature Method	Variance Threshold
SFS RR	Feature Method	Sequential Feature Selection - Ridge Regression
SFS KNN	Feature Method	Sequential Feature Selection - K-Nearest Neighbors
FI RF	Feature Method	Feature Importance - Random Forest
PFI RR	Feature Method	Permutation Feature Importance - Ridge
PFI RF	Feature Method	Permutation Feature Importance - Random Forest
Elastic Loss	Feature Method	ElasticNet Logistic Regression Loss
Elastic Expected	Feature Method	ElasticNet Logistic Regression Expected
Elastic Gain	Feature Method	ElasticNet Logistic Regression Gain
Lasso Loss	Feature Method	Logistic Regression L1 (Lasso) Loss
Lasso Expected	Feature Method	Logistic Regression L1 (Lasso) Expected
Lasso Gain	Feature Method	Logistic Regression L1 (Lasso) Gain
Ridge Loss	Feature Method	Logistic Regression L2 (Ridge) Loss
Ridge Expected	Feature Method	Logistic Regression L2 (Ridge) Expected
Ridge Gain	Feature Method	Logistic Regression L2 (Ridge) Gain

Permutation Feature Importance (PFI) is a technique that replaces the values of a feature with noise and measures the change in performance metrics (such as accuracy) between the baseline and permuted data set. This method overcomes some limitations of impurity-based feature importance but is biased by the correlation between features Hooker et al. (2021). Our ultimate feature set comprises features

exhibiting positive mean importance as determined by the PFI, identifying crucial features. We utilize Random Forests **PFI RF** and Logistic Regression with Ridge Regularization **PFI RR**, both of which assign non-zero scores to all features.

Recursive Feature Elimination (RFE) is a method of training a model on a complete set of features in the data set, eliminating the features with the smallest coefficients. This process iterates until the 10-fold cross-validation score of the models with Random Forest **RFE RF** and Logistic Regression with Ridge Regularization **RFE RR** on the training data shows a decrease. The final scores are attribute rankings where 1 indicates the most relevant features Abe (2005).

Logistic Regression with Filtering and Regularization is a technique that uses L1 **Lasso** or L1 and L2 **Elastic** penalty terms to shrink the coefficients during training. L1 regularization reduces the coefficients of some features to zero for both, and the remaining non-zero coefficients are considered useful information for prediction. On the other hand, L2 regularization, commonly known as **Ridge**, penalizes the square of coefficients, effectively reducing their magnitude without necessarily setting them to zero. This method helps handle multicollinearity and stabilize the model by smoothing out fluctuations in the data, thereby improving generalization performance.

Feature Importance Random Forest (FI RF) is a method that leverages the Random Forests machine learning algorithm to determine the importance of each feature. This importance is measured using either the Gini or the mean decrease impurity. The selected set contains features with the top 50% scores.

Variance Threshold (VT) is a straightforward method to eliminate features by removing features with low variance in the training data set Ghogh et al. (2019). In this work, the threshold used is $0.8 \times (1 - 0.8)$, meaning that features with 80% similar values in the training data set are not selected. The final set of features consists of the k features with the highest variance. VT, SFS RR, and SFS KNN provide a binary selection of features.

Sequential Feature Selection (SFS) searches for the optimal set of features by greedily evaluating all possible combinations of features. The method works by adding one feature at a time and assessing each subset based on the 5-fold cross-validation score of Logistic Regression with Ridge Regression **SFS RR** and **SFS KNN** models.

The labels are also used in Figures 2 to 4 to illustrate the Section and comparisons clearly. Figure 3 illustrates the aggregated scoring mechanism detailed in Algorithm 1. This figure emphasizes the innovative approach to combining filter, embedded, and wrapper methods for feature selection, ultimately producing a robust feature importance ranking.

In total, we obtain ten diverse results comprising binary, numerical, and rank scores. We suggest multiple fusion scoring mechanisms for end-users to consider, as detailed in Algorithm 1. First, we look into five approaches that filter out features and rank the features by the binary sum outputs. Next, we take the methods that provide scores for all features and rank the attribute importance based on the sum of absolute scores. We transform the scores into rankings and combine them with the filtering and ranking methods to develop the final feature, which is importance ranking. Figure 3 illustrates the fusion process described above.

Prediction Modeling

In this study, we address whether the public data collected from web sources is sufficient to predict school district learning performance during the COVID-19 years reliably. Thus, we create five basic baseline models: Logistic Regression with Ridge Regularization, Support Vector Machines (SVM), K-Nearest Neighbor (KNN) suitable for nonlinear and non-separable data, Random Forests, and GBDT. Additionally, we explore four advanced GBDT algorithms: XGBoost, LightGBM, CatBoost, and HistGradientBoosting. Since the data aligns with the features of tabular data, we opt for GBDT methodologies due to their demonstrated robustness in handling diverse tabular datasets Shwartz-Ziv and Armon (2022). The gradient-boosted decision tree (GBDT) assembles many weak decision trees and grows them sequentially and iteratively based on the residual modeling from the previous trees.

Input : Feature Selection Importance Scores(binary, numerical)

Output: Final Fusion Importance Ranking

Initialize BinarySumRankings;

Initialize AbsoluteScoreRankings;

foreach *result in Results* **do**

if *result is binary* **then**

 Apply filtering mechanism to extract relevant features;

 Calculate the binary sum output for these features;

 Rank the features based on the binary sum outputs;

 Append the ranked features to BinarySumRankings;

else

 Apply methods to provide scores for all features;

 Calculate the absolute scores for each feature;

 Rank the attribute importance based on the sum of absolute scores;

 Append the ranked attribute importance to AbsoluteScoreRankings;

end

end

Transform the scores from BinarySumRankings and AbsoluteScoreRankings into rankings;

Combine the rankings derived from both methodologies;

Merge the filtering and ranking methods to generate the

 FinalFeatureImportanceRanking;

return *FinalRanking*;

Algorithm 1: Fusion Scoring Algorithm

The GBDT methods handle tricky observations well and are optimized for faster and more efficient fitting using a data sparsity-aware histogram-based algorithm. In contrast to the pointwise split of the traditional GBDT, which is prone to overfitting, the algorithm's approximate gradient creates estimates by creating a histogram for tree splits. As this histogram algorithm does not handle the sparsity of the data, especially for tabular data with missing values and one-hot encoded categorical features, these algorithms improved tree splits. For example, XGBoost uses Sparsity-aware Split

Finding, defining a default direction of tree split in each tree node Chen and Guestrin (2016). The LightGBM provides the Gradient-Based One-Side Sampling technique, which filters data instances with a large gradient to adjust the influence of the sparsity, and Exclusive Feature Bundling combining features with non-zero values to reduce the number of columns Ke et al. (2017). Our ultimate goal is to assess the predictive power of these nine machine-learning models in this real example.

Open Source Data Acquisition and Processing

Open Data Sources

The dataset utilized for this analysis integrates information from nine distinct sources, employing both School District I.D. and County FIPS Code to cover a comprehensive range of 1,165 school districts across 253 counties in Texas. The data frames and their respective sources include: CCD from the National Center for Education Statistics, providing a matrix of 1189 rows by 66 columns; STAAR and TEA from the Texas Education Agency, with dimensions of 1184x217 and 1182x217 respectively; ADA from the Texas Education Agency, comprising 1226 rows by 3 columns; ESSER from the Texas Education Agency, with 1208 rows by 6 columns; Census data from the Census Bureau (2010), with 254 rows by 37 columns; Covid data from USAFacts, featuring 254 rows by 8 columns; LAUS from the U.S. Bureau of Labor Statistics, spanning 254 rows by 13 columns; and additional Covid data from DSHS, providing a matrix of 1216 rows by 7 columns.

State of Texas Assessments of Academic Readiness (STAAR) data was obtained from the Texas Education Agency (TEA) for the fiscal years 2020, 2021, and 2022 for each school district, 2022c (TEA). The STAAR data we collected are the average scores for math and reading tests and the number of students who participated in the grades 3-8 tests. These data also include students' numbers and average scores under various classifications, such as Title 1 participants, economically disadvantaged, free lunch, special education, Hispanic, Black, White, and Asian.

Common Core of Data (CCD), 2022 (NCES) is the primary database on public elementary and secondary education supplied by the National Center for Education Statistics (NCES) in the United States. The CCD provided us with public school characteristics, student demographics by grade, and faculty information at the school district in Texas for the fiscal years 2019, 2020, 2021, and 2022. For example, according to the data acquired, 62.5% of the students attend school in rural areas, 19.8% in town areas, 10.6% in suburban areas, and only 7.1% in the city area.

Texas School COVID-19 campus data comes from the Texas Department of State Health Services (DSHS), 2022 (DSHS), including the self-reported student enrollment and on-campus enrollment numbers of the dates September 28, 2020, October 30, 2020, and January 29, 2021, at each school district in Texas **County COVID-19** data on infection and death cases due to Coronavirus for each Texas County was parsed from USAFacts source USAFacts (2022). **The average daily attendance (ADA)** is a sum of attendance counts divided by days of instruction per school district and provided by TEA. **Elementary and Secondary School Emergency Relief (ESSER) Grant** data provided by TEA summarizes COVID-19 federal distribution by TEA

to school districts for the fiscal years 2019, 2020, 2021, and 2022. The **Local Area Unemployment Statistics (LAUS)** data, 2022 (BLS) was parsed from the U.S. Bureau of Labor Statistics (BLS) for the years 2019 and 2021 to examine the workforce impact on learning loss in the counties. **Census block group 2010** data Bureau (2010) captures the county's general population characteristics. Upon completing the initial data integration process, merging data from nine sources by matching school district I.D. and county FIPS code, the dataset encompasses 1,165 school districts in Texas, spanning 253 counties with 506 features, one definite, and 505 numerical variables.

For the academic years 2018-2019, 2020-2021, and 2021-2022, the percentage distribution of students by race varied slightly over time. In 2018-2019, the breakdown was as follows: 1.07% Asian or Asian/Pacific Islander, 6.50% Black or African American, 40.41% Hispanic, and 49.09% White. The following year, 2020-2021, saw minor shifts with 1.10% Asian or Asian/Pacific Islander, 6.35% Black or African American, 41.32% Hispanic, and 48.12% White. For the academic year 2021-2022, the proportions were 1.12% Asian or Asian/Pacific Islander, 6.29% Black or African American, 41.67% Hispanic, and 47.70% White.

CARES ESSER I 20, ARP ESSER III 21 features are part of the Elementary and Secondary School Emergency Relief (ESSER) grant programs, which are federal funds granted to State education agencies (SEAs) providing Local education agencies (LEAs) to address the impact due to COVID-19 on elementary and secondary schools across the nation; thus, the funds have been administered by Texas Education Agency (TEA) and allocated in each school district in Texas, 2022 (ESE);, 2021 (TEA). **CARES ESSER I:** Authorized on March 27, 2020, as the Coronavirus Aid Relief and Economic Security (CARES) Act with \$13.2 billion for the fiscal year 2020. **CRRSA ESSER II:** Authorized on December 27, 2020, as the Coronavirus Response and Relief Supplemental Appropriations (CRRSA) Act with \$54.3 billion for the fiscal year 2021. **ARP ESSER III:** Authorized on March 11, 2021, as the American Rescue Plan (ARP) Act with \$122 billion for the fiscal year 2021. **ESSER-SUPP:** Authorized by the Texas Legislature to provide additional resources for not reimbursed costs to support students not performing well educationally from March 13, 2020, to September 30, 2022. To help policymakers make more informative decisions on learning recovery with localized efforts in each school district, we collected data from nine different sources to determine the qualitative conclusions from small sample datasets. Qualitative findings on the educational impacts of COVID-19 highlight significant disruptions in learning environments and the varied responses of educational systems worldwide. It emphasizes qualitative insights into the socio-emotional and pedagogical challenges faced by students and educators during the pandemic Donnelly and Patrinos (2022), matching the data-driven findings from significant, heterogeneous, noisy public data sources.

Data Cleaning, Aggregation and Filtering

Common Core of Data (CCD) for Education Statistics (NCES) is the primary database on public elementary and secondary education supplied by the National Center for Education Statistics (NCES) in the United States. The CCD provided us with public school characteristics, student demographics by grade, and faculty information at the

Table 2. Demographic Proportions of Race/Ethnicity, Gender, and Age Groups Across Counties Bureau (2010).

Race		Gender		Age (0-24)		Age (25+)	
Category	Total	Category	Total	Category	Total	Category	Total
White	17,701,487	Male	12,472,234	0-4	1,928,470	35-44	3,458,373
Black	2,979,598	Female	12,673,245	5-9	1,928,232	45-54	3,435,322
Asian	964,596			10-14	1,881,881	55-64	2,597,668
Hispanic	9,460,903			15-19	1,883,121	65-74	1,472,248
				20-24	1,817,069	75-Up	1,129,626
				25-34	3,613,469		

school district in Texas for the fiscal years 2019 and 2021. Then, we merged the CCD data with the State of Texas Assessments of Academic Readiness (STAAR) data (TEA) from the Texas Education Agency (TEA) for the fiscal years 2019, 2020, 2021, and 2022 at each school district. The STAAR data we collected are the average scores for math and reading tests and the number of students who participated in the grades 3-8 trials. These data also include students' numbers and average scores under various classifications, such as Title 1 participants, economically disadvantaged, free lunch, special education, Hispanic, Black, White, and Asian. Next, our data merged with COVID-19 campus data from the Texas Department of State Health Services (DSHS) of State Health Services (DSHS), including the self-reported student enrollment and on-campus enrollment numbers of the dates September 28, 2020, October 30, 2020, and January 29, 2021, at each school district in Texas. Additional COVID-19 data involved confirmed infection and death cases USAFacts (2022) due to Coronavirus at each county from USAFacts. Also, the average daily attendance (ADA), 2022a (TEA), which consists of the sum of attendance counts divided by days of instruction, and data from the Elementary and Secondary School Emergency Relief (ESSER) Grant Programs (TEA) were collected from TEA for school district level. The ADA data for fiscal years 2019 and 2021 capture the impact of district attendance, and the ESSER data reflect the localized efforts of TEA allocating the grant amount at each school district for the fiscal years 2019, 2020, 2021, and 2022. Also, we combined the Local Area Unemployment Statistics (LAUS) data of Labor Statistics (BLS) from the U.S. Bureau of Labor Statistics (BLS) for the years 2019 and 2021 to examine the negative impact of the unemployment rate on learning loss at the county level. The census block group 2010 data Bureau (2010) capture the demographic features of a county for the general population. We merged the data from nine sources by matching the school district I.D. and county FIPS code and then integrated based on the district I.D. and county FIPS code.

The analysis aggregates the values from various demographic and educational categories into one consolidated group by calculating the percentage difference between the corresponding pairs of values from 2022-2021, 2020-2021, and 2018-2019. For instance, the percentage difference between the total count of White Students in 2020-2021 and 2018-2019, denoted as '% White Students Diff,' reflects the change in the demographic composition over the two years. Similarly, the percentage difference in total race/ethnicity counts and school-wide Title I program participation provides insights into broader demographic shifts and changes in program

enrollment. Moreover, examining the percentage difference in enrollment counts across different grade levels, such as Prek and Kinder, between the two years offers valuable information regarding enrollment trends within specific age groups National Institute for Early Education Research (2021).

Among the 506 features analyzed, 416 display missing values across three data sources, varying from one to 88% within our dataset. Notably, 408 features originate from STAAR and TEA data, six from CCD and NCES, and two from COVID and DSHS data. Within these 416 features, 332 have less than 20% missing values, while 24 exhibits more than 80% missing values. Features exhibiting over 20% missing values primarily originate from the STAAR data, specifically concerning average scores and participation in the STAAR tests. Consequently, we eliminated these features from the STAAR dataset. Additionally, we excluded school districts lacking CCD and COVID data, resulting in 955 public school districts in Texas available for analysis, featuring a total of 119 features devoid of missing values.

Out of 119 features, we aggregate the 58 features that duplicate the data for 2019, 2021, and 2022 into 29 differential features. For example, Total Schools 2020-2021 and Total Schools 2018-2019 features are aggregated into Total Schools Diff 2021, reducing the total number of features to 90. For Total Schools 2022-2021 and Total Schools 2020-2021, features are aggregated into Total Schools Diff 2022, reducing the total number of features to 90. The dataset's missing value distribution across columns reveals varying frequencies within distinct count ranges.

Data Labeling

Our data set is unlabeled; thus, the process begins by normalizing the individual grade scores, ensuring consistency across different scales, through the equation: Normalized Score = (grade score)/max(grade score). Next, the district average is calculated by summing up the scores of grades G3 to G8 and dividing by the total number of grades. The normalized score provides an overarching view of the academic performance within the district, represented by the equation District Average = (G3 + G4 + G5 + G6 + G7 + G8)/(Total number of grades). The percentage loss in performance is computed over time intervals, reflecting changes in educational outcomes, summarized as the score computation in Eq 3.

$$\text{Score} = \frac{\text{Avg 2021} - \text{Avg 2019}}{\text{Avg 2019}} \quad \text{Score} = \frac{\text{Avg 2022} - \text{Avg 2021}}{\text{Avg 2021}} \quad (3)$$

We label the scores as follows: **Learning Gain** if the **Score** is more significant than zero, **Expected** if the **Score** equals zero, and **Learning Loss** if the **Score** is less than zero. This comprehensive process enables the assessment of educational trends, facilitating informed decisions and interventions to enhance learning outcomes.

When analyzed by year, the normalization process encompasses various facets of educational institutions, such as the count of operational public schools, identification of School-wide Title 1 designations, and Title 1 eligibility. The data provide insights into the educational workforce, encompassing Full-Time Equivalent (FTE) teachers and overall staff counts, along with lunch program statistics like free and reduced-price lunch participants. Race and ethnicity distributions among Asian, Hispanic,

Black, and White demographics, delineated by grade groups from Prekindergarten to Grade 12, are normalized for accurate assessment. Attendance metrics undergo normalization in terms of average daily attendance (ADA) and as a percentage of total students per district. By grade, the standardization involves the Percentage of students taking the STAAR reading and math tests, with average scores ratio-ed to the 100th percentile in each grade, regarding population metrics, normalization factors in confirmed COVID-19 cases, and deaths as percentages of the county population. It also encompasses race/ethnicity and age group distributions as a percentage of the county population in 2010. Lastly, when assessed by date, the normalization process considers the Percentage of students on campus on September 28, 2020, October 30, 2020, and January 29, 2021. The Census block grouped by county Bureau (2010) categorizes different household types and housing units as percentages of the total number of households and housing units in 2010, respectively. This comprehensive standardization methodology ensures a consistent and comparable analysis across diverse data points and time frames. Figure 2 illustrates the race, grade, and age groups across counties, highlighting diversity and composition, which are crucial for understanding educational demographics and planning educational resources.

For math, the average learning gain from 2021 to 2022 was 2.80%, contrasting with the previous average loss of -2.75%. This shift indicates an overall improvement in math proficiency. The standard deviation increased to 11.93%, suggesting more significant variability in student outcomes. The minimum loss observed was -50.10%, while the maximum gain was 210.09%. In reading, the average learning gain from 2021 to 2022 was 0.58%, slightly higher than the previous period's gain of 0.32%. The standard deviation remained similar at 9.08%, showing consistent variability. The minimum observed loss in reading was -50.48%, and the maximum gain was 191.43%. Considering all subjects combined, the average learning gain from 2021 to 2022 was 1.69%, a significant improvement compared to the previous average loss of -1.22%. The standard deviation increased to 10.22%, indicating more diverse student outcomes. The minimum observed loss across all subjects was -36.35%, while the maximum gain was 193.44%. Figure 5 shows that student learning outcomes improved from 2021 to 2022, with average gains recorded across all subjects. Math showed the most substantial recovery, transitioning from an average loss to a notable gain, while reading maintained a modest improvement. The increased standard deviations indicate more varied student experiences during this period. Next, we set the threshold to categorize districts based on the STAAR scores into three categories: Loss, Expected, and Gain. The data revealed that more districts experienced a loss in math, with a median loss value of -0.03, compared to a median of 0 for reading. We analyzed and predicted outcomes for math and reading separately. School districts in the middle 50% of loss values were labeled as **Expected**, those <25% as **Loss**, and those 75% as **Gain**.

This categorization enabled us to explore the correlation between various features and these labels. Our findings indicated that the proportion of White students was higher in districts labeled as *Gain* and decreased in those labeled as **Loss**. Conversely, Hispanic students constituted about two-thirds of the **Loss** category, and their proportion decreased in the "Expected" and **Gain** categories for both math and reading. The locale of school districts showed a correlation with learning loss labels as over half of the schools were located in rural areas, with rural locales positively

correlated with the **Gain** label. However, an increasing number of losses occurred in schools located in city and suburban areas.

Data Pre-Processing

In **LossA**, we propose a dimensionality reduction dataset to enhance interpretability and pinpoint resilience factors associated with Learning loss. For instance, features such as "Total Schools 2020-2021" and "Total Schools 2018-2019" are combined into a single feature, "Total Schools Diff," resulting in a total reduction to 90 features. In the dataset **LossB**, these features are treated independently. In this approach, the raw integrated data is employed for the GBDT experiment without normalization while also considering missing values. **LossB** treats each feature individually, whereas **LossA** utilizes normalized and aggregated features to reduce dimensionality and enhance interpretability. While this approach **LossA** may result in a more prominent feature space and potentially increase computational complexity, it allows for a more detailed analysis of features and their impact on learning loss. By examining each feature in isolation, we aim to gain insights into the specific factors contributing to learning loss without the influence of normalization or aggregation techniques. **LossB** provides a complementary perspective to **LossA** and allows for a comprehensive exploration of the dataset, which encompasses 506 features across 1,165 school districts.

Results

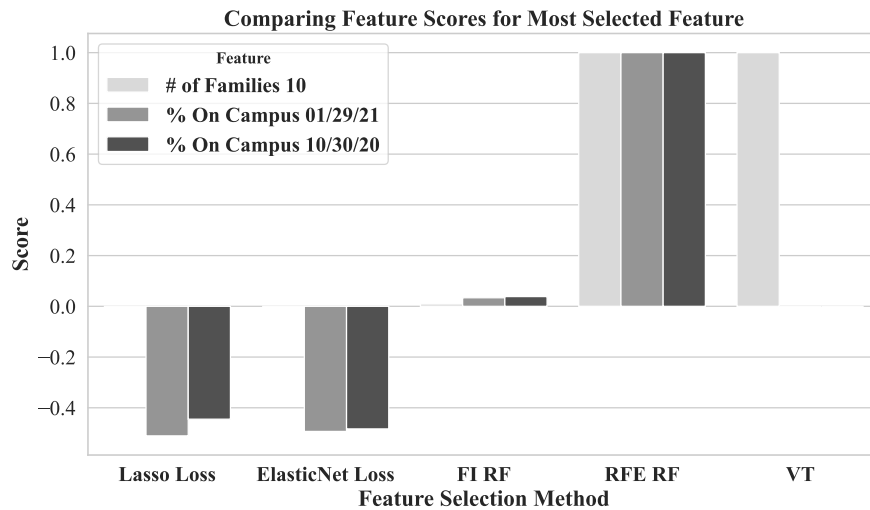


Figure 1. Most Selected Features: # of Families 10, % of Campus 10/30/20, and % On Campus 01/29/21.

In this section, we analyze the results and proposed approaches. For simplicity, Table 1 describes the ten abbreviation labels used for the feature importance scoring.

RFE (RF and RR) provide attribute ranking, and SFS (KNN and RR) provide a binary selection of features. RF FI and PMI (RF and RR) offer non-zero scores to all 90 features. Lasso, Ridge, and Elastic fit for the **Gain**, **Expected**, and **Loss** provides scores for a subset of coefficients selected.

Note that we consider Math Learning Loss and Reading Learning Loss to be two separate tasks with separate attribute selections from the same dataset. Table 1 expands on the following abbreviated feature selection methods that separately detect the resilience factors; abbreviations can be found in Table 1, for Learning **Loss** due to COVID-19 using the data set with 90 features and 955 school districts in Texas as a baseline.

Attribute Importance Analysis

Following Algorithm 1 and the process illustrated in Figure 3, we aggregate five filtering method outcomes for reading and math: VT, SFS KNN, SFS Ridge, and Elastic Gain and Elastic Loss binarized coefficients. The Initial Importance Values represent raw scores from machine learning methods, which are challenging to compare due to their non-uniformity. The Binary Selection Values are the first output transformation, where VT, SFS KNN, and SFS Ridge outcomes are already binary. RFE methods retain the rank of one feature, assigning it a value of 1, and logistic regression assigns +1 to features with positive coefficients and -1 to those with negative coefficients, ignoring coefficients of 0. Feature importance selects the top 50% of features with positive scores as 1, and permutation feature importance assigns 1 to features with positive scores and 0 otherwise. Summing these scores and sorting creates feature importance rankings for each subject out of 9.

The Impact Score Values, the second transformation, normalize scores by dividing each method's scores by the sum of all feature scores to ensure equal contribution to the final ranking. The summation of the absolute value of normalized scores forms feature ranking. The top 20 features with the highest scores are selected for math and reading, prioritizing the impact score, which integrates binary and non-zero scores. The binary score serves as a secondary measure for understanding importance, determining the cutoff point where the impact score drops. Secondary labels are applied to features to categorize their type, enhancing understanding of their significance. This approach facilitates the comparison of feature importance and identification of the most critical features in educational data analysis.

In this study, we applied several feature selection methods to determine their effectiveness in reducing the dimensionality of features for math and reading scores. The methods evaluated include RFE (RF and RR), VT, SFS (RF and KNN), FI RF, PFI (RR and RF), Elastic, Lasso, and Ridge, (*Loss*, *Gain*, *Expected*). Lasso, Elastic, PFI RR, PFI RF, and FI RF generate score outputs. Lasso resulted in 51 features for both math and reading. Figure 6 shows methods SFS KNN, SFS RR, and VT that focus on selecting subsets of features based on iterative addition or variance criteria. SFS methods iteratively add features to improve model performance, while VT removes low-variance features to enhance efficiency and accuracy. Figure 4 illustrates RFE RF and RFE RR employ recursive techniques to systematically reduce feature sets, facilitating the identification of the most influential features

while optimizing computational efficiency. Grouping these methods aids in visualizing the interpretability of selected features by each feature selection method. Figures 2 include F1 RF, PFI RR, and PFI RF, which emphasize feature importance assessment through recursive elimination or permutation testing. These methods evaluate how features contribute to model accuracy, identifying crucial predictors for refined model performance.

Figure 1 illustrated the most selected features such as # of *Families 10*, % *On Campus 01/29/21*, and % *On Campus 10/30/20* by feature selection techniques listed in Table 1). The # of *Families 10* showed significant relevance across all methods, indicating its robust influence on educational outcomes during disrupted learning environments. The %*On Campus 01/29/21* and % *On Campus 10/30/20* demonstrated varying impacts depending on the method employed.

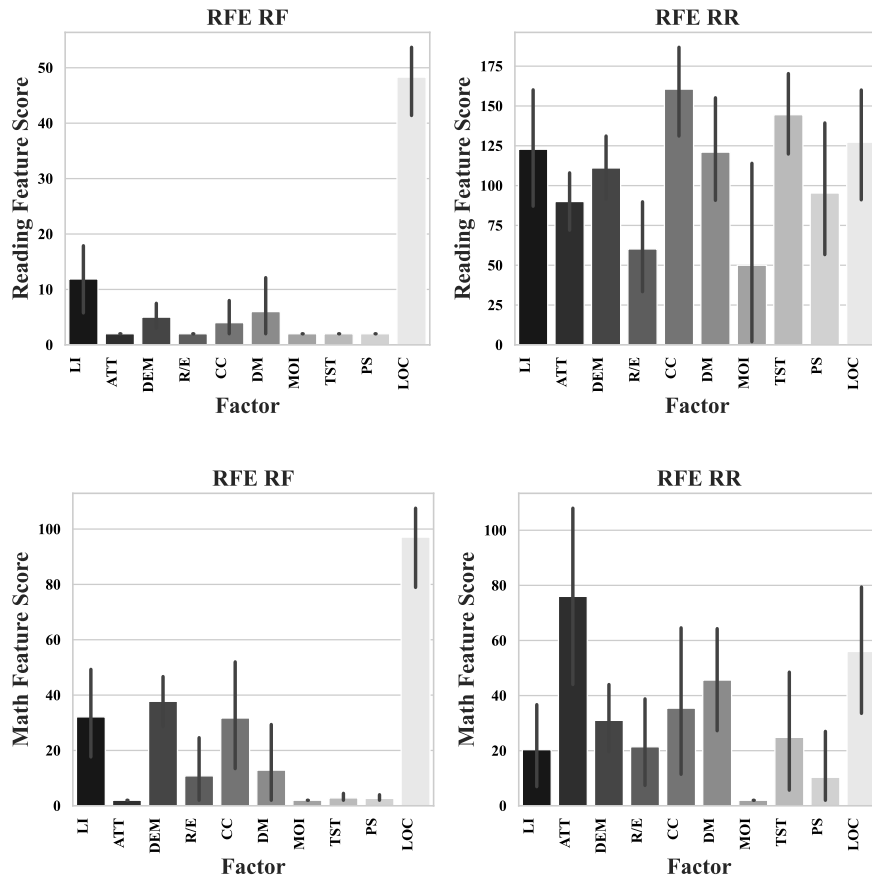


Figure 2. Filtering Feature Selection Methods for Reading (top) and Math (bottom) Comparison of RFE RF and RR.

The ten methods ranked 18 features as top importance and agreed to exclude 33 descriptors, mainly from the workforce, Census, and COVID data sources. The difference between free lunch and the COVID deaths in the county had little impact on learning loss. Next, we sort the remaining 57 features using RF FI, PMI (RF and RR), RFE (RR and RF) scores, and Elastic Gain and Elastic Loss. Since all of them have an importance ranking per feature (including the sign), we first normalize the scores for each method and then sum them as listed in Table 1.

Elastic reduced the features to 45 for reading and 41 for math. PFI RR identified 82 for reading and 28 features for math, while PFI RF selected 26 for reading and 70 features for math. FI RF resulted in 45 features for both subjects. Methods producing binary outputs include VT, SFS RR, and SFS KNN. VT reduced the features to 20 for both reading and math. Both SFS RR and SFS KNN selected 45 features for each subject. Methods producing rank outputs include RFE RR and RFE RF. RFE RR resulted in the most miniature feature set, with five features for reading and six for math, whereas RFE RF identified 36 features for both subjects. The Permutation Feature Importance (PFI) methods identified the most significant number of features, with PFI RF selecting 70 features for math and PFI RR selecting 82 features for reading. Table 6 presents the importance ranking of the features and summarizes the top 20 features for math (a) and for reading (b) selected by six or more methods scores in 2021 and 2022. The most significant feature predicting learning loss in math is *% of Campus 10/30/20*, the enrollment of students in the campus district on October 30, 2020, representing the mode of instruction.

For reading, three critical features were selected, all of which were resilience factors related to the Low-income backgrounds of students: *CARES ESSER I 20* (Coronavirus Aid, Relief and Economic Security (CARES) grant amount in 2020), *ARP ESSER III 21* (American Rescue Plan Act (ARP) grant amount in 2021), *% Reduced-price Lunch Diff* (Reduced-price Lunch Eligible Students Difference in percent between 2019 and 2021). Table 5 summarizes the top 20 features of impact to the learning loss. Table 6 summarizes the top 20 features for learning recovery. The attribute is important if selected by six or more selection methods summarized in Table 1. Figure 4 illustrate that income and Grade level are the most influential resilient factors to predict learning loss for math and reading. The race/Ethnicity and mode of instruction continued to be decisive, resilient factors for both subjects; on the other hand, Attendance and Census demographics are considered significant factors only in math, and Unemployment is essential only for reading. Although we now realize these primary features can identify the resilient factors for **Loss** or Gain in learning due to the COVID-19 pandemic, it is still unknown whether those features positively impact learning. We analyzed positive or negative correlations between the most critical features and our label, Loss, Expected, or Gain, in math and reading. The students who experienced **Loss** in reading received more significant funding for all funding programs on average than the students who participated, showed Gain or Expected in the same subject. The districts in need of financial help for adapting and preparing for learning **Loss** due to COVID-19 received the ESSER funds amounts calculated by a formula based on Title I and Part A grant (ESE, TEA).

Figure 4 indicates that *% of Campus 10/30/20* is positively correlated with Gain as the Distribution of school districts with the highest proportion of students on a campus

populated more for Gain and Expected in math; however, the students experienced **Loss** are inhabited the most where the enrollment is 0%. It is clear that in-person classes, the mode of instruction, were the key to avoiding **Loss** in math.

*Modeling Learning **Loss** from Public Data*

The data sets have been randomly split into 80% of the training set and 20% of the test set with shuffling and stratification on the label. We use performance metrics suitable for prediction problems to find the best model. The accuracy score for both Gain and **Loss** is used to get a big picture, and the F1 score is used for an in-depth measure as it harmonically includes the precision and the recall scores. Matthews correlation coefficient (MCC) considers true negatives, class imbalance, and multi-class data. Each model runs with a 10-fold cross-validation of GridSearch to find optimal hyperparameters. The boosting algorithm trains weak learners iteratively, and early stopping reduces training time and avoids overfitting. At every boost round, the model evaluates and decides whether to stop or continue the training when it shows no more improvement for a certain number of consecutive rounds in terms of the evaluation metric specified as the fit parameter. The number of early stopping rounds is set to 10% of the maximum number of boosting iterations.

We employed five state-of-the-art machine learning models: Ridge Regression, Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Random Forests, and GBDT and their variations, as summarized in Table 1. We trained the models using our complete set of 90 features and ten additional feature groups derived from various feature selection techniques, and Table 1 summarizes the model characteristics. Table 4 outlines the performance of the machine learning methods. Performance metrics, including accuracy, F1 score, and Matthews correlation coefficient (MCC), for these models are presented in bar graphs in Figure 3 for baseline models and in Figure 4 for GBDT models. The prediction of learning loss for reading exhibits weaker performance than for math. Although most models perform similarly across both subjects, with the exception of KNN, GBDT for math, and ridge regression for reading, they demonstrate the highest average accuracy, F1 score, and MCC.

We penalize and regularize the algorithm by hyperparameter tuning so that we aim to increase accuracy and avoid overfitting to improve the gradient boosting modeling for XGBoost, LightGBM, CatBoost, and HistGradientBoosting. These hyperparameters are searched with a 5-fold cross-validation RandomizedSearch with the number of iterations that is 20% of parameter distributions of each model. XGBoost is supposed to explore 100 distributions of the parameters; the number of iterations for RandomizedSearch is 20 times. The constraints on tree structures aid in curbing the growth of overly complex and extensive trees, limiting the number of trees, tree depth, and the number of leaves per tree in the model. A lower learning rate (below 0.5) allows for a gradual adjustment of tree weights during each iteration, thereby minimizing errors. Ridge and Lasso regularization terms further the models by simplifying the complexity and size of the model Chen and Guestrin (2016). The GBDT algorithms also show higher prediction power for math than reading and indicate no significant model exceeding other models, including the best state-of-the-art models, in terms of performance.

Table 3. Five State-of-the-Art Machine Learning Models Evaluated on the test set (20% of the data)

Model	Accuracy	Precision	Recall	F1	MCC	ROC
GB	0.6329	0.6131	0.6329	0.5877	0.3513	0.7339
KNN	0.5911	0.5731	0.5911	0.5548	0.2698	0.6631
RF	0.6339	0.6355	0.6339	0.5651	0.3602	0.7288
Ridge	0.6273	0.5974	0.6273	0.5662	0.3453	0.7270
SVM	0.6268	0.5982	0.6268	0.5588	0.3441	0.7257

For math, the best-performing model was CatBoost, achieving an accuracy of 67.5%, an F1 score of 64.5%, and an MCC of 43.4% using 36 features selected by RFE RF. The other notable performances included Gradient Boost with an accuracy of 64.4%, an F1 score of 62.2%, and an MCC of 37.5% using the same feature selection method, and XGBoost with an accuracy of 66.0%, an F1 score of 61.6%, and an MCC of 40.5% using 21 features selected by Variance Threshold (VT).

For reading, CatBoost also emerged as the top performer, achieving an accuracy of 62.3%, an F1 score of 54.8%, and an MCC of 33.8% using 82 features selected by PMI Ridge. The second-best performance was from XGBoost, with an accuracy of 61.3%, an F1 score of 53.5%, and an MCC of 31.2% using 90 features from all feature sets, followed by Ridge, which had an accuracy of 60.7%, an F1 score of 52.2%, and an MCC of 30.3% using 45 features selected by SFS Ridge.

Overall, the GBDT algorithms, CatBoost and XGBoost, were the best choices among all the machine learning models tested for predicting learning loss in both subjects. Despite better performance in predicting math scores than reading scores, the performance gap between the four GBDT models and the five other state-of-the-art models, except KNN, was negligible, with a difference in accuracy of around 3%.

Additionally, no single dimensionality reduction technique consistently outperformed others across all models. The various dimensions of the selected features were experimented with to examine the effects of dimensionality reduction methods and the best set of the features by predicting learning loss with the machine learning models introduced in Section .

Best features vs. Raw Data for GBDT Models

In this section, we analyze the performance of four GBDT models—XGBoost, LightGBM, CatBoost, and HistGradientBoosting—on different datasets to evaluate their predictive power regarding learning loss in math and reading. These models inherently handle data sparsity, including missing values, by finding optimal tree splits. The initial dataset, referred to as **LossB**, consists of 506 features (505 numerical and one categorical) across 1,165 school districts, with 416 features containing missing values ranging from 1% to 88%.

We compared the performance of models trained on three datasets: (1) the best feature sets identified through various feature engineering techniques, (2) raw data without imputation for missing values, and (3) raw data with missing values imputed using mean values. The subject-specific features differ for math and reading. Each subject ended up with 302 features, with 212 features containing missing values.

Table 4. Four Advanced GBDT Models (and their st.dev.) evaluated on the test set

Model	Accuracy	Precision	Recall	F1	MCC
CatBoost	0.6337 (0.03)	0.6310 (0.07)	0.6337 (0.03)	0.5778 (0.05)	0.3543 (0.07)
HistGB	0.6304 (0.03)	0.6157 (0.06)	0.6304 (0.03)	0.5805 (0.04)	0.3464 (0.07)
LightGBM	0.6281 (0.03)	0.6117 (0.06)	0.6281 (0.03)	0.5735 (0.04)	0.3415 (0.07)
XGBoost	0.6247 (0.03)	0.6047 (0.07)	0.6247 (0.03)	0.5635 (0.04)	0.3363 (0.06)

The comparison in Table 4 showed that all models improved their performance metrics, especially the MCC when using the best feature sets compared to raw data. *istGradientBoosting* exhibited the most significant improvement in MCC for math, increasing by 47%, followed by *LightGBM* 43%, *CatBoost* 25%, and *XGBoost* 24%. In reading, the improvement in MCC was even more pronounced, with *HistGradientBoosting* showing a 124% increase and *LightGBM*, *CatBoost*, and *XGBoost* improving by 45%, 43%, and 41%, respectively. Additionally, models trained on raw data without imputation performed slightly better than those with imputed data across all subjects and models. MCC for math increased by over 6% for *CatBoost* and *HistGradientBoosting*, while *XGBoost* showed the most significant MCC growth for reading, with an increase of 10%. In conclusion, gradient-boosted decision tree (GBDT) models trained on carefully selected feature sets significantly outperformed those trained on raw data, highlighting the importance of feature engineering in predictive modeling. Moreover, avoiding the imputation of missing values yielded better performance than mean imputation, emphasizing the models' capability to handle raw data effectively. Table 4 also illustrates that over ten feature selection methods, the GBDT models are robust against changes in feature subsets as the standard deviation of the results (in brackets is usually 1%). The models maintain similar levels of performance regardless of the specific features used for training, which is beneficial in ensuring reliable predictions across different datasets or real-world applications.

Conclusion and Future Work

In this study, we employ a data-driven approach to investigate the impact of the COVID-19 pandemic on learning loss, utilizing an intentional data science pipeline. Despite employing ten distinct feature selection methods to facilitate the automatic extraction of crucial features from publicly available datasets, our findings reveal a limited influence on prediction accuracy across the nine machine learning models trained on both feature-selected sets and raw data. Notably, GBDT algorithms, particularly *XGBoost* and *CatBoost*, consistently outperform other models, demonstrating remarkable efficacy in managing missing values prevalent within the raw datasets. Your reproducible experiments and datasets are accessible at Yu and Tešić (2022), providing valuable tools for policymakers to strategically allocate resources and interventions to mitigate the effects of learning loss. A deeper analysis of 2021 to 2022 data revealed that shifts in feature significance primarily occurred at the individual feature level rather than through changes in resilience factor importance. Significantly, the mode of instruction and prior score emerged as the primary resilience factors during this period. Overall, low income and grade level proved to be the most influential

factors in predicting learning loss in both math and reading. Noteworthy contributors to math performance include attendance and census demographics, particularly the % of Campus 10/30/20. Additionally, students from low-income backgrounds and regions with higher unemployment rates were particularly impactful in predicting reading learning loss. In future research, we aim to broaden the temporal scope of our analysis and incorporate more granular data sources to deepen our understanding of the enduring repercussions of the COVID-19 pandemic on education. Additionally, exploring novel feature engineering techniques or enhancing existing ones could bolster prediction accuracy across various datasets.

References

- Abe, S. (2005). Modified backward feature selection by cross validation. In *ESANN*, pages 163–168, U.S. Springer.
- Abutbul, A., Elidan, G., Katzir, L., and El-Yaniv, R. (2020). Dnf-net: A neural architecture for tabular data. *CoRR*, abs/2006.06465.
- Agency, T. E. (2021). Impacts of covid-19 and accountability updates for 2022 and beyond. <https://tea.texas.gov/sites/default/files/2021-tac-accountability-presentation-final.pdf>.
- Arik, S. and Pfister, T. (2021). Tabnet: Attentive interpretable tabular learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(8):6679–6687.
- Baashar, Y., Alkawsi, G., Ali, N., Alhussian, H., and Bahbouh, H. (2021). Predicting student’s performance using machine learning methods: A systematic literature review. In *2021 International Conference on Computer & Information Sciences (ICCOINS)*, pages 357–362, U.S. IEEE.
- Betebenner, D., Van Iwaarden, A., Cooperman, A., Boyer, M., and Dadey, N. (2021). Assessing the academic impact of covid-19 in summer 2021.
- Bureau, C. (2010). Census block group 2010. <https://schoolsdata2-93b5c-tea-texas.opendata.arcgis.com/datasets/census-block-group-2010-tx/>.
- Cardona, T., Cudney, E. A., Hoerl, R., and Snyder, J. (2020). Data mining and machine learning retention models in higher education. *Journal of College Student Retention: Research, Theory & Practice*, 25(1):1521025120964920.
- Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’16, page 785–794, New York, NY, USA. Association for Computing Machinery.
- Choate, K., Goldhaber, D., and Theobald, R. (2021). The effects of covid-19 on teacher preparation. *Phi Delta Kappan*, 102(7):52–57.
- Courtemanche, C. J., Le, A. H., Yelowitz, A., and Zimmer, R. (2021). School reopenings, mobility, and covid-19 spread: Evidence from texas. Technical report, National Bureau of Economic Research.
- Donnelly, R. and Patrinos, H. A. (2022). Learning loss during covid-19: An early systematic review. *Prospects*, 51:601–609.
- (ESE) (2022). Elementary and secondary school emergency relief fund. <https://oese.ed.gov/offices/education-stabilization-fund/>

- elementary-secondary-school-emergency-relief-fund/.
for Education Statistics (NCES), N. C. (2022). Common core of data (ccd).
<https://nces.ed.gov/ccd/elsi/tableGenerator.aspx>.
- Ghojogh, B., Samad, M., Mashhadi, S., Kapoor, T., Ali, W., Karray, F., and Crowley, M. (2019). Feature selection and feature extraction in pattern analysis: A literature review. *arXiv:1905.02845v1*, 1(02845):1–14.
- Grinsztajn, L., Oyallon, E., and Varoquaux, G. (2022). Why do tree-based models still outperform deep learning on tabular data?
- Halloran, C., Jack, R., Okun, J. C., and Oster, E. (2021). Pandemic schooling mode and student test scores: Evidence from us states. Technical report, National Bureau of Economic Research.
- Hooker, G., Mentch, L., and Zhou, S. (2021). Unrestricted permutation forces extrapolation: variable importance requires at least one more model, or there is no free variable importance. *Statistics and Computing*, 31:1–16.
- Joseph, M. (2021). Pytorch tabular: A framework for deep learning with tabular data.
- Ke, G., Meng, Q., Finely, T., Wang, T., Chen, W., Ma, W., Ye, Q., and Liu, T.-Y. (2017). Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems 30 (NIP 2017)*, pages 1–9, <https://www.microsoft.com/en-us/research/publication/lightgbm-a-highly-efficient-gradient-boosting-decision-tree/>. Advances in Neural Information Processing Systems 30 (NIP 2017).
- Ke, G., Zhang, J., Xu, Z., Bian, J., and Liu, T.-Y. (2019). TabNN: A universal neural network solution for tabular data.
- Kim, Y., Choi, Y.-K., and Emery, S. (2013). Logistic regression with multiple random effects: A simulation study of estimation methods and statistical packages. *The American Statistician*, 67(3):171–182.
- Maldonado, J. E. and De Witte, K. (2022). The effect of school closures on standardised student test outcomes. *British Educational Research Journal*, 48(1):49–94.
- National Institute for Early Education Research (2021). Impacts of COVID-19 on Preschool enrollment and spending.
- OECD (2021). *Education at a Glance 2021*. Organisation for Economic Co-operation and Development, <https://doi.org/10.1787/b35a14e5-en>.
- of Labor Statistics (BLS), U. B. (2022). Local area unemployment statistics (laus). <https://www.bls.gov/lau>.
- of Public Instruction, N. C. D. (2022).
- of State Health Services (DSHS), T. D. (2022). Texas public schools covid-19 data. <https://dshs.texas.gov/coronavirus/schools/texas-education-agency/>.
- Popov, S., Morozov, S., and Babenko, A. (2019). Neural oblivious decision ensembles for deep learning on tabular data. *CoRR*, abs/1909.06312:1–12.
- Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., and Gulin, A. (2018). Catboost: unbiased boosting with categorical features. *Advances in neural information processing systems*, 31.
- Rao, A. R., Desai, Y., and Mishra, K. (2019). Data science education through education data: an end-to-end perspective. In *2019 IEEE Integrated STEM Education Conference (ISEC)*, pages 300–307, U.S. IEEE.

- Shwartz-Ziv, R. and Armon, A. (2022). Tabular data: Deep learning is not all you need. *Information Fusion*, 81:84–90.
- (TEA), T. E. A. (2021). Elementary and secondary school emergency relief (esser) grant programs. <https://tea.texas.gov/finance-and-grants/grants/elementary-and-secondary-school-emergency-relief-esser-grant-programs>.
- (TEA), T. E. A. (2022a). Average daily attendance (ada). <https://tea.texas.gov/finance-and-grants/state-funding/state-funding-reports-and-data/average-daily-attendance-and-wealth-per-average-daily-attendance>.
- (TEA), T. E. A. (2022b). Impacts of covid-19 and accountability updates for 2022 and beyond.
- (TEA), T. E. A. (2022c). State of texas assessments of academic readiness (staar) for 2018-2019 and 2020-2021. <https://tea.texas.gov/student-assessment/testing/staar/staar-aggregate-data>.
- USAFacts (2022). Texas coronavirus cases and deaths. <https://usafacts.org/visualizations/coronavirus-covid-19-spread-map/state/texas>.
- Wang, Z. (2020). When large-scale assessments meet data science: The big-fish-little-pond effect in fourth-and eighth-grade mathematics across nations. *Frontiers in Psychology*, 11:579545.
- Yan, K. (2021). Student performance prediction using xgboost method from a macro perspective. In *2021 2nd International Conference on Computing and Data Science (CDS)*, pages 453–459, U.S. IEEE.
- Yu, J. and Tešić, J. (2022). Tabular data in the wild: Gradient boosting modeling improvement. <https://github.com/DataLab12/educationDataScience>.
- Zamarro, G., Camp, A., Fuchsman, D., and McGee, J. B. (2022). Understanding how covid-19 has changed teachers’ chances of remaining in the classroom. *Sinquefield Center for Applied Economic Research Working Paper*, 22(01).

Supplemental material

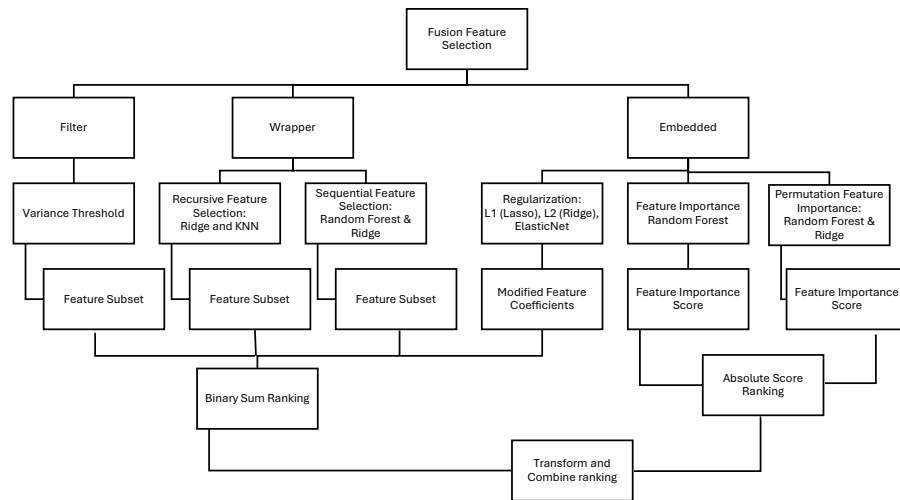


Figure 3. Fusion Process Of Aggregating Outcomes From Various Feature Selection Methods For Reading And Math For The Academic Years 2018-2019, 2020-2021, 2021-2022. *Mirna TODO: unreadable, if you cant fix the source send it to me to fix.*

Table 5. Top 20 Math Features for 2021 and 2022 Datasets

Math Scores 2021			Math Scores 2022		
Feature	Impact Score	Binary Score	Feature	Impact Score	Binary Score
Median Household Income	6.62	5	Average Annual Pay	6.40	3
Total Students 2018-2019	6.23	7	Per Capita Income	6.27	4
Total Students 2020-2021	6.14	6	Total Students 2021-2022	6.02	6
Total Students 2021-2022	6.11	7	County Population	5.92	5
Rural: Distant	6.05	3	# of Families 10	5.91	6
# of Families 10	5.84	4	Total Students 2018-2019	5.89	5
Average Annual Pay	5.83	2	Total Students 2020-2021	5.87	5
ARP ESSER III 21 NORM	5.76	3	# of Households 10	5.84	5
CARES ESSER I 20 NORM	5.76	4	% of Population Under 18 in Poverty	5.80	4
Rural: Remote	5.74	3	CRRSA ESSER II 21 NORM	5.81	4
# of Housing Units 10	5.70	3	Median Household Income	5.78	5
# of Households 10	5.70	3	# of Housing Units 10	5.78	4
Per Capita Income	5.70	3	Median Age Female 10	5.76	3
% of Population Under 18 in Poverty	5.68	3	% of Population in Poverty	5.77	4
Median Age Male 10	5.68	3	Rural: Distant	5.70	3
County Population	5.68	2	CARES ESSER I 20 NORM	5.71	4
% of Population in Poverty	5.67	2	ARP ESSER III 21 NORM	5.69	4
CRRSA ESSER II 21 NORM	5.67	2	Median Age Male 10	5.66	3
Median Age 10	5.65	2	Median Age 10	5.59	2
Median Age Female 10	5.58	1	Rural: Remote	5.56	2

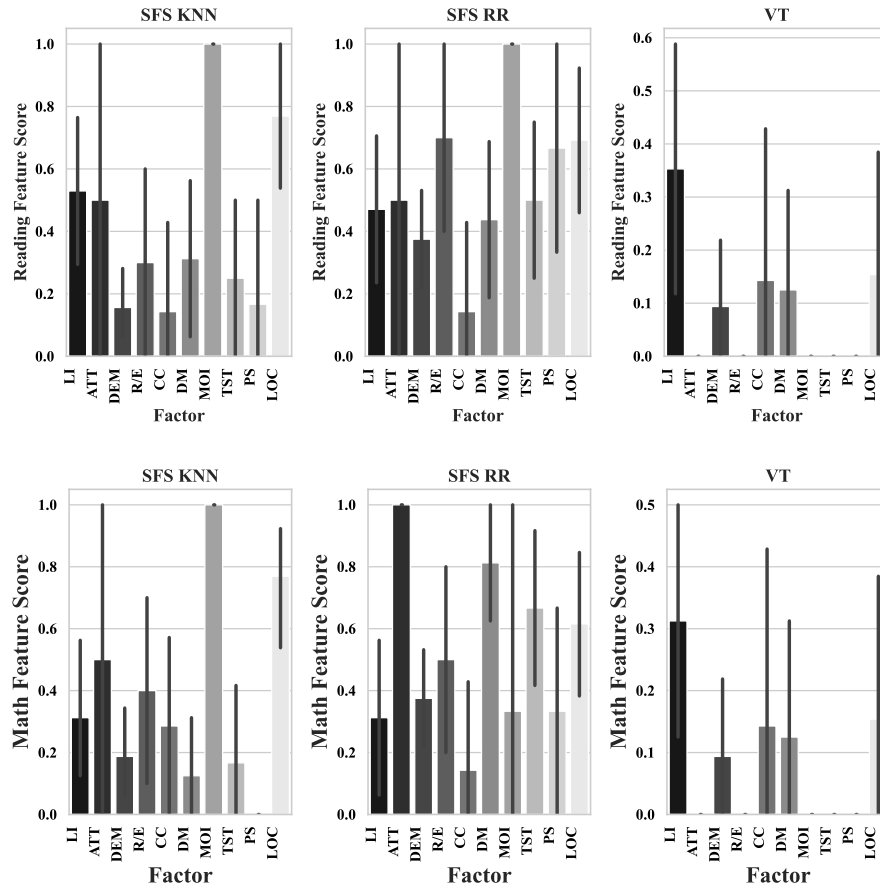


Figure 4. Importance Feature Selection Methods for Reading (top) and Math (bottom): Comparison of SFS KNN, SFS RR, and VT for identifying significant resilience factors.

Table 6. Top 20 Reading Features for 2021 and 2022 Datasets

Reading Scores 2021			Reading Scores 2022		
Feature	Impact Score	Binary Score	Feature	Impact Score	Binary Score
Median Household Income	5.78	5	Total Students 2018-2019	5.89	4
Total Students 2018-2019	5.89	4	Total Students 2020-2021	5.87	5
Total Students 2020-2021	5.87	5	# of Households 10	5.84	5
Total Students 2021-2022	5.85	5	% of Population Under 18 in Poverty	5.80	4
# of Households 10	5.84	5	CRRSA ESSER II 21 NORM	5.81	4
# of Housing Units 10	5.78	4	Median Household Income	5.78	5
Median Age Female 10	5.76	3	# of Housing Units 10	5.78	4
% of Population in Poverty	5.77	4	Median Age Female 10	5.76	3
Rural: Distant	5.70	3	% of Population in Poverty	5.77	4
CARES ESSER I 20 NORM	5.71	4	Rural: Distant	5.70	3
ARP ESSER III 21 NORM	5.69	4	CARES ESSER I 20 NORM	5.71	4
Median Age Male 10	5.66	3	ARP ESSER III 21 NORM	5.69	4
Median Age 10	5.59	2	Median Age Male 10	5.66	3
Rural: Remote	5.56	2	Median Age 10	5.59	2

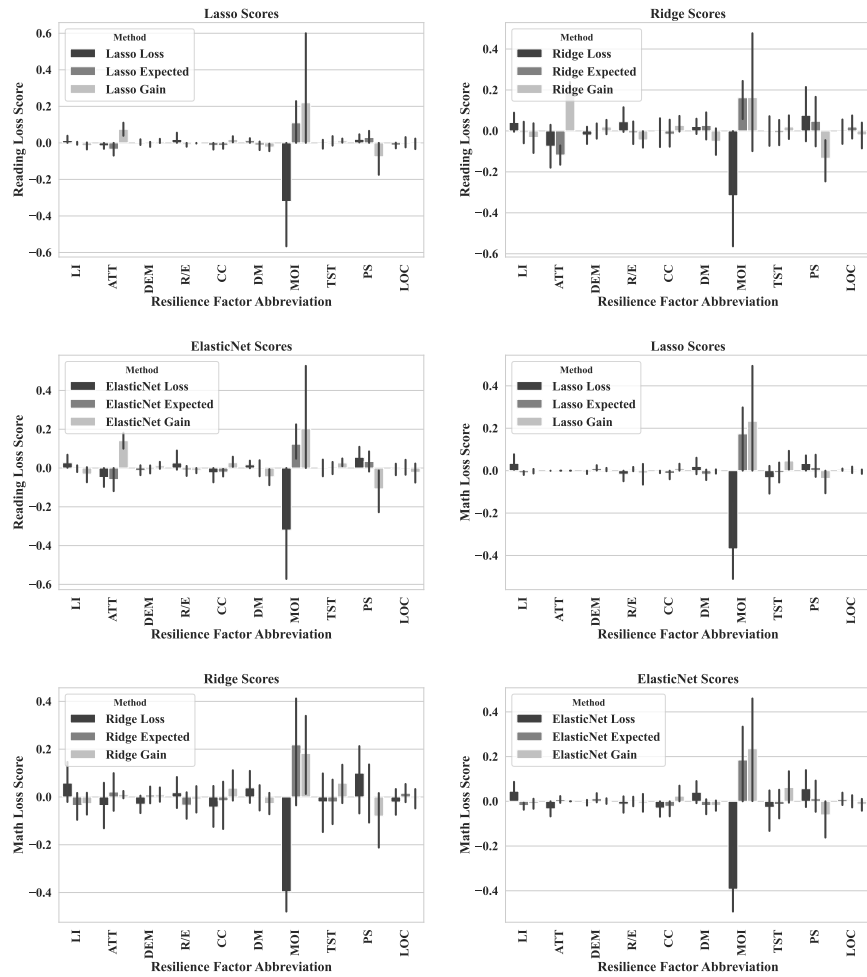
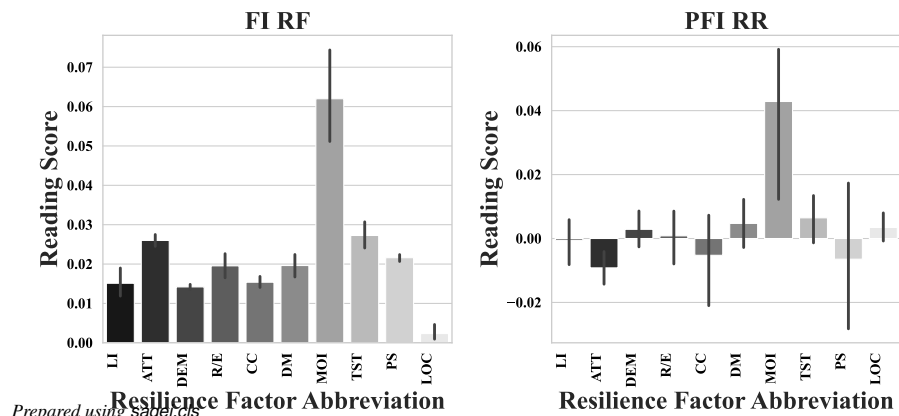


Figure 5. Comparison of Loss, Expected, and Gain Metrics for Lasso, Ridge, and Elastic Regularization Techniques to identify significant reading (top) and math (bottom) resilience factors



Prepared using sagej.cls