

1

2

3

4

5

Applying XGB Regression Trees to Produce Growth

6

Percentiles

7

8

Steven Tang, Zhen Li

9

eMetric LLC

10

11

12

13

14

Paper written for the 2019 meeting of the National Council on Measurement in

15

Education, Toronto, Canada. The views expressed in this paper are solely those of the

16

authors and they do not necessarily reflect the positions of eMetric LLC.

17

Correspondence concerning this paper should be addressed to Steven Tang, eMetric,

18

211 N Loop 1604 E, Suite 170, TX 78232. Email: steven@emetric.net.

19 **Abstract**

20 This study compares percentile rank residuals using an XGBoost regression tree model
21 to quantile regression based SGP. Results indicate that with default hyperparameters,
22 the XGB tree based approach can exactly replicate standard SGP, and that the XGB
23 method may be further tuned to potentially predict more accurately.

24 *Keywords:* Gradient boosted regression tree, growth percentile ranking, student
25 growth percentile

26 **Background**

27 In recent years, big data methods such as gradient boosted decision trees and
28 deep neural network architectures have shown great promise in tackling a variety of
29 prediction modeling tasks, often surpassing the results from traditional methods or
30 even solving previously unsolvable prediction tasks. In this study, we investigate the
31 potential for applying gradient boosted regression trees, enabled through the XGB
32 statistical package, to the prediction task of computing student growth measures, under
33 the hypothesis that the favorable statistical properties of XGB models may allow for an
34 alternative procedure to compute growth measures similar to quantile regression based
35 student growth percentiles (SGP).

36 SGP has been used for measuring students' annual growth in many states. In
37 theory, an SGP describes a student's relative progress with respect to his/her academic

38 peers, who are students beginning at the same place (Betebenner, 2008, 2018). Quantile
39 regression is commonly used to estimate the conditional growth percentiles of current-
40 year scores based on prior year scores.

41 Castellano & Ho (2013) explored using percentile rank residuals (PRR) based off
42 of ordinary least squares (OLS) regression and found that the OLS regression method
43 proved to be a promising alternative to the quantile-regression based SGP method, as
44 the OLS regression PRR method recovers the true conditional status percentile ranks
45 better in certain situations. However, OLS regression is known to have strict
46 assumptions such as homoskedasticity of the errors and gaussian distributions of the
47 covariates, et al.

48 In this study, eXtreme Gradient Boosting (XGB) regression trees using the PRR
49 method are applied to two case study datasets as an alternative to the quantile-
50 regression based SGP approach. Both XGB and quantile regression relax the
51 homoskedasticity assumption, but XGB goes a step further and makes no assumption
52 that data distributions need to be gaussian or that relationships must be linear.
53 Moreover, XGB regression trees have favorable properties such as high predictive
54 accuracies with many possible input variables, a tweakable and tunable training
55 procedure, fast computation, and an interpretable decision-tree structure that can be
56 illuminated after training. XGB approaches can be prone to issues of overfitting,
57 therefore requiring special consideration in model construction and interpretation.

76 calculated by XGB PRR and SGP for an extensive comparison. Table 2 presents the
 77 descriptive statistics of the second dataset.

78 Table 2

79 *Descriptive Statistics of 2018 Mathematical and ELA Test Data*

Cohort	Year	Grade	Mathematics			ELA		
			N	Mean	S.D.	N	Mean	S.D.
1	2017	3	37803	2426	81	37868	2418	84
	2018	4	38311	2465	81	38309	2467	85
2	2016	3	37626	2423	79	37682	2420	83
	2017	4	38089	2463	81	38099	2461	87
	2018	5	38684	2489	89	38776	2498	90
3	2016	4	36241	2459	79	36306	2462	86
	2017	5	36804	2488	85	36868	2499	90
	2018	6	37459	2500	99	37527	2512	89
4	2016	5	35475	2485	84	35518	2499	85
	2017	6	36105	2497	98	36147	2511	89
	2018	7	36528	2509	106	36585	2539	96
5	2016	6	34631	2498	97	34684	2508	84
	2017	7	34654	2506	101	35361	2539	95
	2018	8	35502	2524	111	35577	2555	98

80

81 **XGB Regression Trees**

82 The XGB Regression Tree approach relies on iteratively building a collection of
 83 simple regression trees; regression trees are decision trees that predict continuous
 84 outcomes. The iterative process starts by first creating an extremely simple predictive
 85 regression tree; such a tree might only have between 2 to 16 leaf nodes. This initial
 86 regression tree is constructed by searching through a large number of potential split

87 values among all input variables and finding the splits that minimize prediction error.
88 The iterative process continues by constructing an additional regression tree of the same
89 structure, but this time constructed to minimize the *residual errors* of the first regression
90 tree. The next iterative tree is then constructed to minimize the residuals of the full
91 model thus far, and the process of iteratively creating new trees continues until
92 stopping criteria is met. As the name implies, gradient boosting uses gradient descent to
93 find the next regression tree to add to the ensemble. At the end of the building process,
94 the predictions are given by the sum of the outputs of all trees. This process of building
95 a gradient boosted regression tree was optimized in the XGB package allowing for very
96 fast computation of gradient boosted trees as well as many opportunities for additional
97 model tuning (Benjamin, Fernandes, Tomlinson, Ramkumar, VerSteeg, Miller, &
98 Kording, 2014).

99 For a predictive model $\hat{y}_1 = f_1(X)$, where X indicates input variables, \hat{y}_1
100 indicates predictions by the first tree and y indicates the observed output variable, a
101 loss function can be defined between the prediction and the observed outcome: $l(\hat{y}_1, y)$.
102 During training, the first tree can be estimated by minimizing the following objective:

$$L_1 = \sum l(\hat{y}_1, y) + \Omega(f_1) \quad (1)$$

103 Ω is a regularizing function to avoid overfitting. Then a second tree $f_2(X)$ will be
104 constructed by predicting the residuals of the first tree. The objective to minimize is as
105 follows:

$$L_2 = \sum l(\hat{y}_1 + f_2(X), y) + \Omega(f_2) \quad (2)$$

106 The process continued sequentially for a fixed number of trees (N). Total loss will be
107 progressively decreased with each additional tree. In the end, the prediction for y will
108 be the sum of the predictions of all trees:

$$\hat{y} = \sum_k^N f_k(X) \quad (3)$$

109 Compared to linear regression and quantile regression, XGB regression tree
110 require completely different assumptions. For example, linear regression has a basic
111 assumption that the sum of its residuals is 0. XGB regression tree, through its boosting
112 process, instead attempts to find and model patterns in the residuals and strengthen the
113 model with weak learners that exploit these patterns. This approach has shown to be
114 extremely powerful in big data tasks, winning a variety of competitions where
115 predictions need to be made based on a wide set of predictors.

116 **Procedure of Applying XGB to Produce Percentile Ranks of Residual**

117 To produce XGB PRR, the following steps were carried out: 1) Train a XGB
118 prediction model with two or more years of consecutive scale scores for one cohort of

119 students; 2) Use the prediction model to generate a predicted score, which is regarded
120 as the expected score that a student should have got in the current year; 3) Compute a
121 current-year residual score by subtracting the predicted score from the current-year
122 observed score; 4) Calculate PRR, the percentage of students whose residual scores are
123 lower than or equal to the score of interest in the population. A function “rankdata”
124 from a python package “scipy.stats.mstats” is used to compute ranks (order statistics) of
125 each residual score. When the residual scores are tied, the average rank is used. Then
126 the following formula is applied to compute percentile ranks.

$$PRR = \text{round}\left(100 \times \frac{\text{rank}_x - 1}{N}\right) \quad (4)$$

127 Equation (4) is slightly different from the equation (4) in Castellano & Ho’s (2013)
128 article, where they calculated PRR as the percentage of residual scores that are smaller
129 or equal to the score of interest. Another definition of percentile rank is the percentage
130 of residual scores less than the target score plus 0.5 of the percentages of ties in all
131 residual scores. The different definitions of percentile ranks might lead to slightly
132 different outcomes, but these differences should be minor after we round the
133 percentages to integers. In addition, PRR is forced to be located within [1,99] to compare
134 to SGP.

135 The XGB results presented in this study use the XGB package (Chen & Guestrin,
136 2016) implemented in Python. SGP results are obtained using the SGP package
137 (Betebenner, 2018) in R. Results from two studies are presented in the following section.

138 Results and Discussion

139 The first result comes from comparing XGB PRR and SGP using just two years of
140 scale scores for a state mathematical test. In **Error! Reference source not found.**, four
141 different models' results are shown, each trained to incorporate different input
142 variables. A hyperparameter grid search was performed to mitigate overfitting
143 concerns. Results show that the base model, where only grade 7 math is used to predict
144 grade 8 math scores, can achieve a 0.997 correlation with standard SGP.

145 However, as more input variables are incorporated, the correlation with SGP
146 goes down, but R^2 with realized scores correspondingly increases. This means that the
147 XGB PRR model with more input variables disagrees more with SGP, but has better
148 model predictive accuracy relative to realized scores. This provides evidence that it is
149 operationally easy for the XGB PRR approach to replicate standard (quantile regression)
150 SGP results, but that incorporating additional explanatory variables can increase model
151 accuracy and correspondingly decrease correlation with standard SGP.

152 A trained XGB regression tree model can also be inspected to better understand
153 how the model is making decisions. There are numerous metrics that can be used.

154 Figure 1 depicts the most important features used in the most complex model, which
 155 used grade 7 math, grade 7 reading, and demographic variables as predictors.

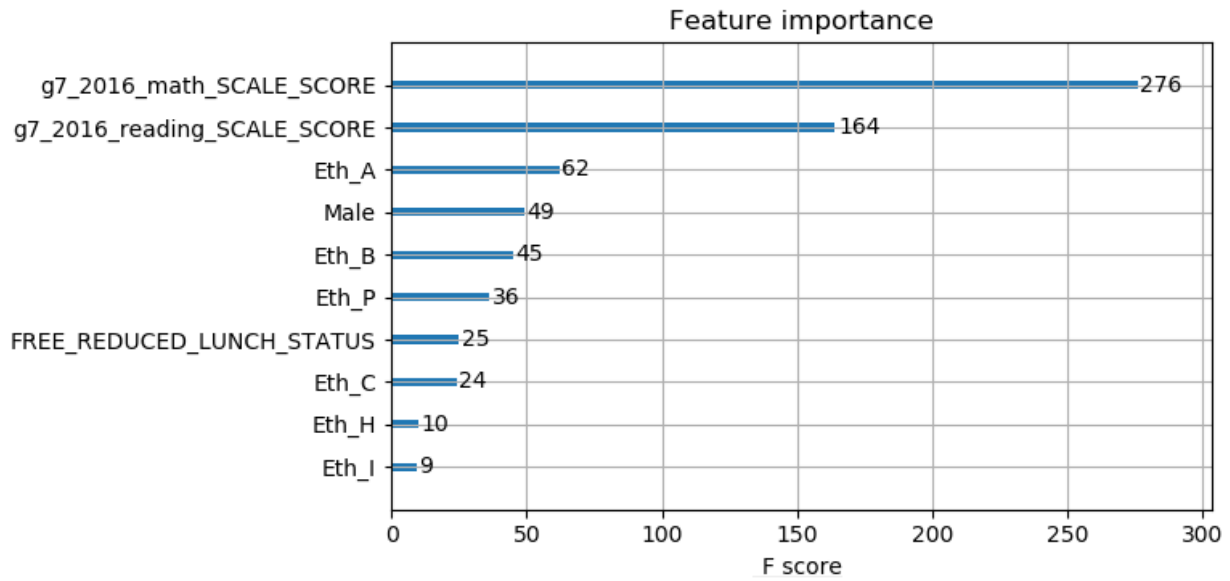
156 Previous studies (Castellano & Ho, 2013; Lockwood & Castellano, 2015) found
 157 that alternative estimation methods (OLS based SGP or Logit model based SGP) can
 158 provide SGP estimates closer to the SGP calculated using empirical conditional
 159 distributional functions (ECDF). We didn't use ECDF in the current study, although this
 160 may be useful to look at in future studies. Lockwood and Castellano (2015) also showed
 161 that even if the correlation between the estimates by different methods are very high,
 162 the small difference between individual SGP estimates can cause significant effect for
 163 teacher evaluation, which is based on group-level SGP.

164 Table 3

165 *XGB Model Results in the Pilot Study (2016-2017 Mathematical Test)*

Input to XGB Model	Hyperparameters (All but the first model was chosen via 5-fold Cross-Validation)	Correlation with SGP	R ²
G7 Math	Estimators = 100, Max Depth = 1, Learning Rate = .1	.997	.619
G7 Math + Demographics	Estimators = 700, Max Depth = 1, Learning Rate = .04	.985	.628
G7 Math + G7 Reading	Estimators = 600, Max Depth = 1, Learning Rate = .03	.951	.650
G7 Math + G7 Reading + Demographics	Estimators = 700, Max Depth = 1, Learning Rate = .04	.945	.653

166



167

168

Figure 1 XGB Feature Importance from the Pilot Study

169

Next, using 2016-2018 students’ scale score data from both mathematics and ELA

170

test administrations, we compared the two models with more prior years’ scale scores.

171

For XGB PRR, we apply a simple XGB regression tree model with most

172

hyperparameters set as default values. The number of estimators was fixed to 125 and

173

max depth was fixed as 4 for all prediction models.

174

175

176

177

178

179

180

181

182

183 Table 4

184 *XGB Model Results for 2018 Test Data*

Output	Input Variables	Correlation with SGP	R^2
G8 Math	G6 Math+G7 Math	.991	.769
G8 Reading	G6 Reading+G7 Reading	.993	.778
G7 Math	G5 Math+G6 Math	.990	.812
G7 Reading	G5 Reading+G6 Reading	.991	.772
G6 Math	G4 Math+G5 Math	.989	.787
G6 Reading	G4 Reading+G5 Reading	.993	.763
G5 Math	G3 Math+G4 Math	.992	.781
G5 Reading	G3 Reading+G4 Reading	.992	.768
G4 Math	G3 Math	.996	.759
G4 Reading	G3 Reading	.995	.723

185

186 Table 5

187 *XGB Model Results with more Input Variables*

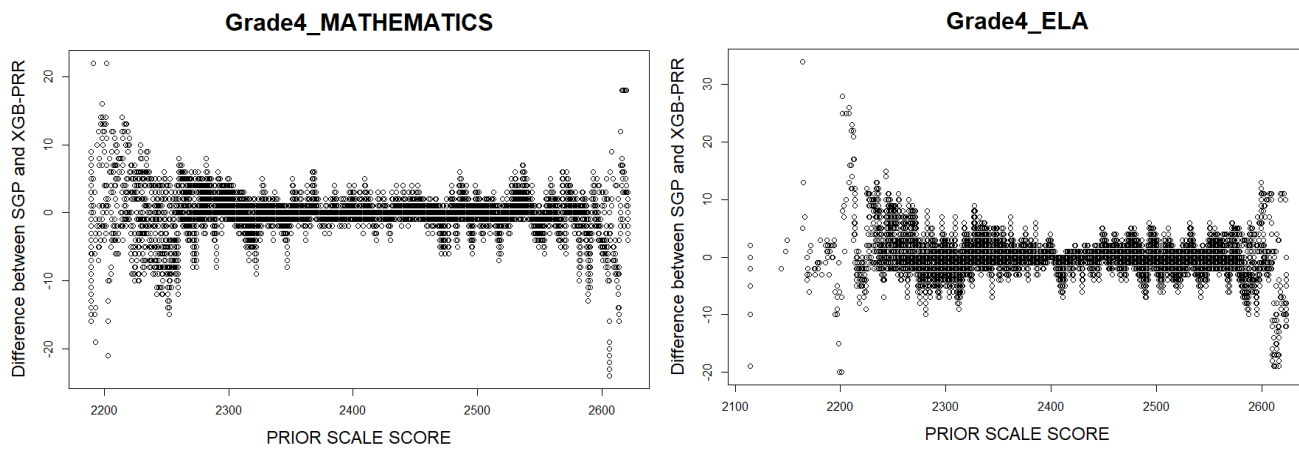
Output	Input Variables	Correlation with SGP	R^2
G8 Math	G6 Math+G7 Math+G6 Reading+G7 Reading + Demographics	.954	.788
G8 Reading	G6 Math+G7 Math+G6 Reading+G7 Reading + Demographics	.957	.794
G7 Math	G5 Math+G6 Math+G5 Reading+G6 Reading + Demographics	.958	.824
G7 Reading	G5 Math+G6 Math+G5 Reading+G6 Reading + Demographics	.956	.788
G6 Math	G4 Math+G5 Math+G4 Reading+G5 Reading + Demographics	.939	.810
G6 Reading	G4 Math+G5 Math+G4 Reading+G5 Reading + Demographics	.960	.779
G5 Math	G3 Math+G4 Math+G3 Reading+G4 Reading + Demographics	.971	.791
G5 Reading	G3 Math+G4 Math+G3 Reading+G4 Reading + Demographics	.958	.784
G4 Math	G3 Math+G3 Reading + Demographics	.966	.775
G4 Reading	G3 Math+G3 Reading + Demographics	.934	.756

188 Results from Table 4 shows that the correlation coefficients between XGB PRR
189 and SGP range from .989 to .996. The correlation coefficients are equivalently high
190 across all grades and subjects. Results in Table 5 shows that when incorporating
191 additional input variables (more subjects and demographics), the correlation between
192 XGB PRR and standard SGP decreased and R^2 increased. These results closely mimic
193 the trend found from Table 3, where adding more data to the XGB model decreased
194 correlation to SGP results but increased overall R^2 .

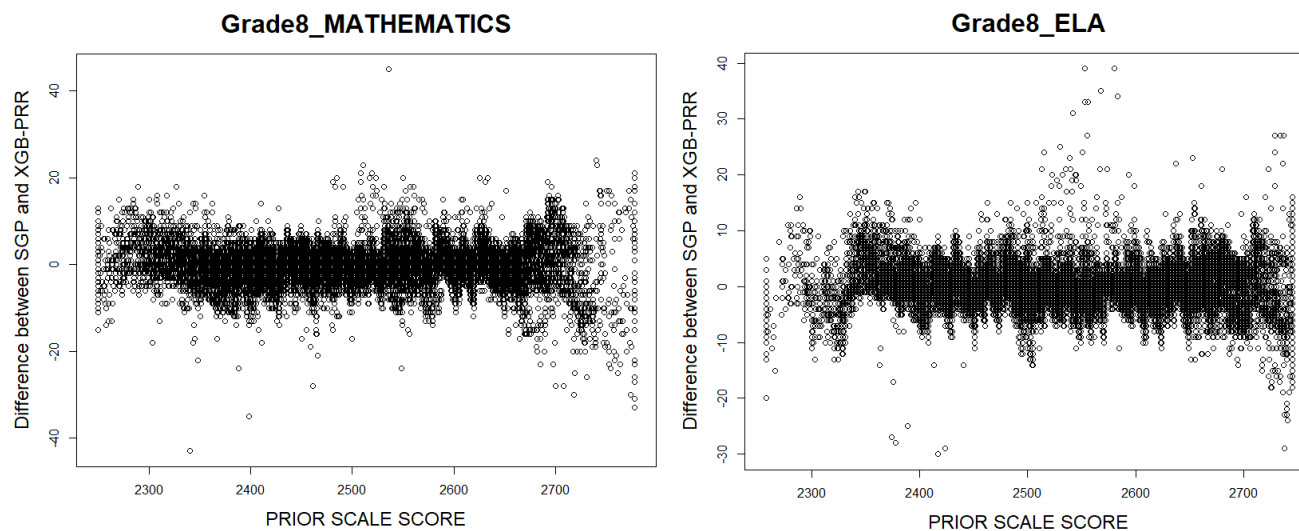
195 Furthermore, results from 2018 data analysis show that the difference between
196 XGB-PRR and SGP are higher at the extreme previous year scale scores. This effect is
197 very significant for Grade 4 tests, where the input variables only include one prior year
198 test data. When the number of prior years increase, this pattern is not as clear. As
199 shown in Figure 2, for grade 8 ELA and math, the largest difference occurs for extreme
200 scoring students, but also shows a little bit in the middle. This effect was also
201 discovered in the 2017 data analysis and in a previous study (Castellano & Ho, 2013).

202

203



204



205

Figure 2. The Difference between two Growth Measures (SGP-XGBPRR) across

206

Prior-year Scale Scores

207

Conclusion

208

The practical purpose of this study is to see if XGB PRR could be a feasible

209

alternative statistical framework to quantile regression SGP. The results of this paper

210

indicate that an XGB based model of ranking student growth can, at a minimum,

211

replicate ranks produced by quantile regression based SGP. However, the XGB model

212 has additional statistical properties that may make it preferable, such as being able to
213 model more input features to achieve better predictive accuracies. Additionally, the
214 XGB framework is easy to operationalize, is robust to missing data, and is relatively
215 easy to interpret and analyze.

216 To establish the XGB PRR as a useful and viable alternative will take additional
217 research, but given how successfully the XGB approach has been applied to many other
218 big data prediction tasks, this line of research appears to be quite promising. There are
219 numerous avenues for future exploration to utilize the expressive and robust properties
220 of the XGB decision tree methodology for prediction. Additionally, other prediction
221 problems in educational statistics, such as making useful forecasts of other results
222 besides growth measures, may also be addressed by modern statistical frameworks like
223 XGB regression trees. The results presented in this study can contribute to a fuller
224 understanding of how modern statistical methods can solve or improve on problems of
225 prediction in large scale measurement.

226

227

228

229

230

231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248

References

- Benjamin, A.S., Fernandes, H.L., Tomlinson, T., Ramkumar, P., VerSteeg, C., Miller, L., & Kording, K.P. (2018). Modern machine learning far outperforms GLMs at predicting spikes. *Frontiers in Computational Neuroscience*, 12 (56), 1-13.
- Betebenner, D. W. (2008). Toward a normative understanding of student growth. In K. E. Ryan & L. A. Shepard (Eds.), *The future of test-based educational accountability* (pp. 155–170). New York, NY: Taylor & Francis.
- Betebenner, D. W. (2018). SGP: Student growth percentile and percentile growth Trajectories [R package version 1.8-0.0].
- Castellano, K. E. & Ho, A. D. (2013). Contrasting OLS and quantile regression approaches to student “growth” percentiles. *Journal of Educational and Behavioral Statistics*, 38(2), 190-215.
- Chen, T. & Guestrin, C. (2016). *XGB: A Scalable Tree Boosting System*. Paper presented in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco.
- Lockwood, J. R. & Castellano, K. E. (2015). Alternative statistical frameworks for student growth percentile estimation, *Statistics and Public Policy*, 2:1, 1-9