

1

2

3

4

5

6

Forecasting students' future academic performance using big data analytics

7

8

9

10

11

Zhen Li, Steven Tang

12

13

eMetric LLC

14

15

16

17

18

19

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

20

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

21

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

22

Correspondence concerning this paper should be addressed to Zhen Li, eMetric, 211 N

23

Loop 1604 E, Suite 170, TX 78232. Email: zli@emetric.net.

24

Abstract

25

26

27

28

29

30

31

Keywords: XGBoost regression tree, Bayesian networks, K-12 assessment

32

33

34

35

36

37

38

39

40

41

42

43

In education, big data analytics methods have become increasingly popular (Romero & Ventura, 2010). This article illustrates how we use XGBoost regression trees for predicting students' future performance in state summative tests. Bayesian networks and linear regression model are applied for comparison. Results show that XGBoost regression trees perform the best, with higher prediction accuracy and computation efficiency. The XGBoost regression tree also works better with incomplete data sets.

Year after year, students take high stakes summative tests, and the results of these tests can have far-reaching consequences for students, teachers, and other stakeholders. In this study, we investigate the possibility of using the XGBoost statistical framework, which implements gradient boosted regression trees, in order to make potentially useful forecasts of student scores on high stakes summative tests. Given the current and prior scores of a particular student, we seek to forecast how that student will do on next year's tests. This type of information could be useful to many stakeholders; teachers and schools could draft a plan to create targeted interventions for at-risk students, for example. The underlying hypothesis is that modern methods such as XGBoost regression have proven to be statistically accurate and operationally easy to use and may be able to provide a feasible statistical framework to provide score

44 forecasts, and such predictions could eventually be disseminated via reporting to
45 various stakeholders. We seek to compare XGBoost results to other commonly used
46 statistical frameworks in education literature, namely Bayesian networks and linear
47 regression. The statistical frameworks will be evaluated using overall predictive
48 accuracy (root-mean -square error) as well as robustness to missing data.

49

50

The Big Data Analytics Models

51 **XGBoost regression tree (XGBoost)**. This approach relies on iteratively building a
52 collection of simple regression trees; regression trees are decision trees that predict
53 continuous outcomes (Chen & Guestrin, 2016). The iterative process starts by first
54 creating an extremely simple predictive regression tree; such a tree might only have
55 between 2 to 16 leaf nodes. This initial regression tree is constructed by searching
56 through a large number of potential split values among all input variables and finding
57 the splits that minimize prediction error. The iterative process continues by constructing
58 an additional regression tree of the same structure, but this time constructed to
59 minimize the *residual errors* of the first regression tree. The next iterative tree is then
60 constructed to minimize the residuals of the full model thus far, and the process of
61 iteratively creating new trees continues until stopping criteria is met. As the name
62 implies, gradient boosting uses gradient descent to find the next regression tree to add
63 to the ensemble. At the end of the building process, the predictions are given by the

64 sum of the outputs of all trees. This process of building a gradient boosted regression
65 tree was optimized in the XGBoost package allowing for very fast computation of
66 gradient boosted trees as well as many opportunities for additional model tuning
67 (Benjamin, Fernandes, Tomlinson, Ramkumar, VerSteeg, Miller, & Kording, 2014).

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

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

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

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

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

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

78 Compared to linear regression and quantile regression, XGBoost regression tree
79 require completely different assumptions. For example, linear regression has a basic

80 assumption that the sum of its residuals is 0. XGBoost regression tree, through its
81 boosting process, instead attempts to find and model patterns in the residuals and
82 strengthen the model with weak learners that exploit these patterns. This approach has
83 shown to be extremely powerful in big data tasks, winning a variety of competitions
84 where predictions need to be made based on a wide set of predictors.

85 **Bayesian networks (BN).** Based upon a joint distribution for a directed acyclic graph,
86 Bayesian networks can estimate conditional probability of one variable given other
87 variables in the net. As we know, building a Bayesian net consists of two parts:
88 structure learning and parameter learning. The structure of a net can be either freely
89 estimated or pre-defined. In this study, we compared results from a learned structure
90 and a fixed structure and found the prediction results very close to each other. With a
91 large number of input variables, structure learning is very time demanding. Therefore,
92 a simple fixed structure was applied for all the Bayesian networks modeling.

$$P(y|\mathbf{X}) = P(y) \prod_{k=1}^n P(x_k|y) \quad (1)$$

93 Where $\mathbf{X} = (x_1, \dots, x_k, \dots, x_n)$ indicates the input variables, y indicates the score field to
94 be predicted. The number of input variables is n . The net only has edges from all the
95 input variables to the target variables, which means that the target variable is
96 dependent on all the input variables. Furthermore, all the input variables are assumed

97 to be independent. The parameters of the structure (conditional probabilities) were
98 freely estimated by maximum likelihood estimation. The R package "bnlearn" is used
99 for parameter calibration (Scutari, 2010). As all functions in "bnlearn" require complete
100 data, the training data only contains students with complete observations. For the test
101 data, we impute the input variables with the learned net at the first step and predict the
102 target variables at the second step.

103 Bayesian networks (Pearl & Scutari, 2000; Scutari, 2010) have been thoroughly
104 studied for several decades and is also popular in the psychometrics field (Pearl &
105 Scutari, 2000; Mislevy, Almond, Yan & Steinberg, 2000; Tsamardinos, Brown, & Aliferis,
106 2006; Sinharay, 2006; Scanagatta, de Campos, Corani, & Zaffalon, 2015). Comparing to
107 other machine learning models, Bayesian networks have shown several advantages.
108 First, expert knowledge of the net structure and conditional probabilities can be
109 incorporated. Second, all the parameters in Bayesian networks are interpretable and can
110 be presented clearly in a graph. Third, no specific input and output variables need to be
111 defined. That is to say, once the net is learned and calibrated, the values of any variable
112 can be predicted using the other variables. Fourth, Bayesian networks have also been
113 found to be robust to missing data (Friedman, 1997). Fifth, likelihoods can be provided
114 to predicted scores. Finally, Bayesian networks have been applied in psychometrics for
115 decades. For example, Mislevy et al. (2000) applied Bayesian networks to model
116 relationships between latent cognitive variables; Sinharay (2006) applied the posterior

117 predictive model checking method to evaluate model fit of Bayesian nets. Therefore, we
 118 select Bayesian networks as our second method.

119 **Methodology**

120 **Data**

121 One cohort of students' test scores in reading, writing, math, and science from
 122 grade 3 to grade 8 were collected. Science was only taken in grade 5 and grade 8. The
 123 following table shows the subjects tested at each grade.

124 Table 1.

125 *Test Data per Grade*

	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Reading	√	√	√	√	√	√
Math	√	√	√	√	√	√
Science			√			√

126 Note: "√" means that the subject was tested at the purported grade.

127 Test scores included scale scores, performance levels, as well as reporting
 128 category scores for each subject. About a quarter of students had incomplete records.
 129 Additionally, students' demographic information, e.g., gender, ethnicity, were also
 130 included in the data input file. In the output variable (predicted field), only valid test
 131 scores were selected. The total number of students in each test ranged from 300,000 to
 132 400,000. 80% of the data was randomly chosen for training and validation, while the
 133 remaining 20% was used as a test dataset.

134 Study Design

135 The aim of this study is to evaluate XGB in predicting students' next-year
136 academic performance in summative tests. We compare XGB with two popular
137 approaches: Bayesian networks and linear regression. In the prediction model
138 framework, the input variables include all previous years' test scores and students'
139 demographic information (2013-2017). The output variables are test scores at the most
140 recent year (2018). For students in a lower grade, e.g., grade 4, only one previous-year
141 data exist (e.g., grade 3 in 2017); However, students in a higher grade, e.g., grade 8, have
142 many more previous years of test data (e.g., grade 3 in 2013 - grade 7 in 2017). In this
143 study, we also explore how the prediction accuracy of XGB could improve when more
144 previous years of test data are used as input variables. In the end, we compare the
145 performance of XGB and Bayesian networks with regard to their prediction accuracy for
146 students with incomplete data.

147 Evaluation Criteria

148 We used root mean squared error (RMSE), mean errors (ME) and classification
149 consistency to evaluate the performance of the prediction models.

$$RMSE = \sqrt{\sum_{i=1}^N (SS_{forecast} - SS_{observed})^2 / N}, \quad (2)$$

$$ME = \sum_{i=1}^N (SS_{forecast} - SS_{observed}) / N, \quad (3)$$

150 where N is the total number of students for a test; $SS_{forecast}$ indicates predicted scale
151 scores; $SS_{observed}$ indicates the observed scale scores.

152 Classification consistency is defined as the probability that the predicted scores
153 and real scores classify students into the same performance level group, based on the
154 given performance level cuts for each test.

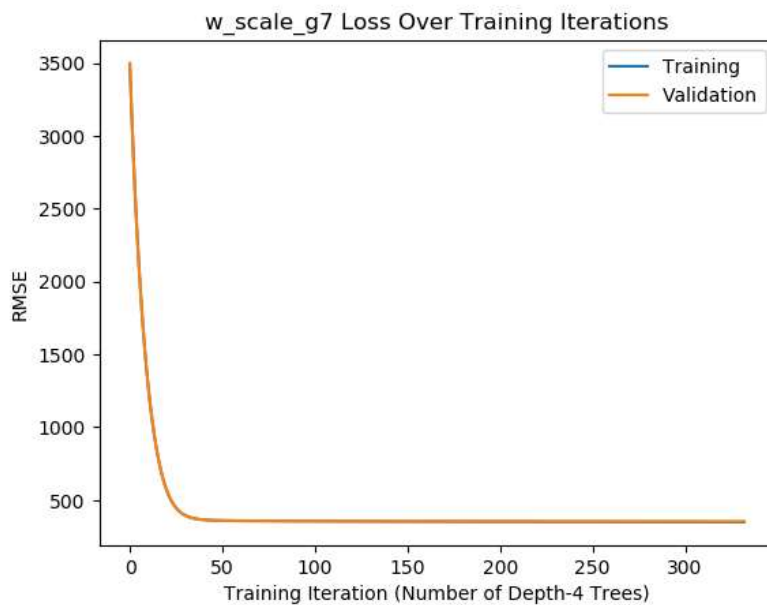
155 **Results**

156 The three above-mentioned methods for predicting students' academic
157 performance were applied to a longitudinal data set, consisting of students' test scores
158 for 6 years in a state assessment. We predicted students' scale scores of different
159 subjects at Grades 4-8 by all their corresponding previous-year data. Results are
160 presented in this section.

161 **Model fit**

162 Psychometric models commonly report one or several model fit indices when applied to
163 real data. However, machine learning packages do not produce model fit indices
164 directly. Usually, machine learning models are evaluated using different training,
165 validation and test datasets. The prediction errors on the validation and test data set are

166 the major criteria of evaluation. XGBoost also produce the loss functions across training.
167 Figure 1 shows an example of the training and validation loss function across iterations
168 by XGBoost regression tree. Prediction errors for the training and validation data
169 decrease at the same time with more iterations, which provides evidence that
170 overfitting doesn't happen. More complex model evaluation, such as cross validation,
171 could also be carried out for both methods. But as our sample size is very large while
172 the number of input variables is relatively small, it is evident that the training,
173 validation and test data in our study are all representative of the full data.



174

175 *Figure 1* Loss over training iterations by XGBoost

176 **Classification Consistency**

177 Using the predicted scores, classification consistency indices were calculated
 178 based on known cut-off scores. From 2012-17, this test has two fixed cut-off standards:
 179 “Performance Level Cut 1” and “Performance Level cut 2”. Table 2 presents the
 180 classification consistency at each performance level cut respectively.

181 Table 2

182 *Comparison of classification consistency index for two performance level cuts*

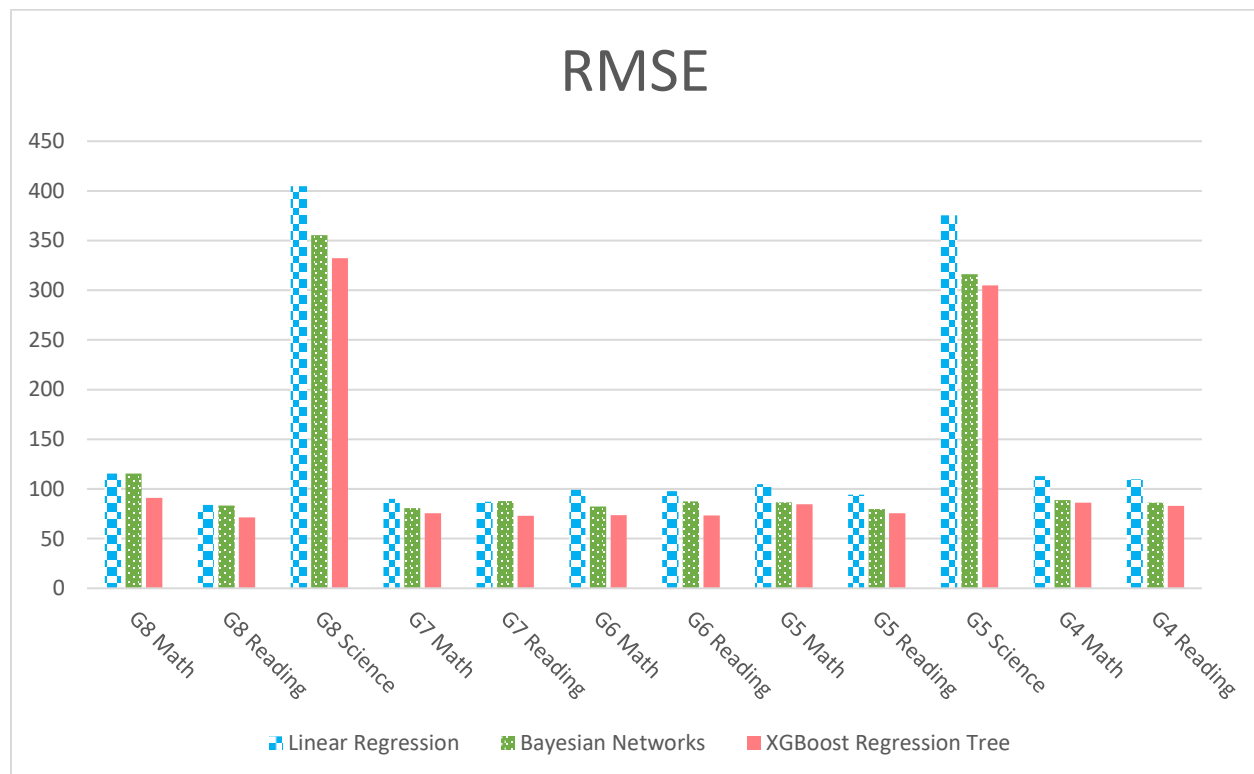
Target Field	Performance Level Cut 1			Performance Level Cut 2		
	Linear regression	Bayesian networks	Regression tree	Linear regression	Bayesian networks	Regression tree
G8 Math	0.698	0.742	0.788	0.882	0.896	0.901
G8 Reading	0.822	0.815	0.845	0.858	0.863	0.874
G8 Science	0.802	0.804	0.818	0.856	0.878	0.885
G7 Math	0.832	0.839	0.853	0.899	0.903	0.909
G7 Reading	0.820	0.819	0.842	0.856	0.866	0.876
G6 Math	0.780	0.831	0.845	0.882	0.910	0.915
G6 Reading	0.786	0.832	0.847	0.846	0.884	0.889
G5 Math	0.784	0.818	0.822	0.863	0.882	0.885
G5 Reading	0.787	0.828	0.833	0.853	0.877	0.880
G5 Science	0.759	0.808	0.810	0.898	0.910	0.911
G4 Math	0.797	0.823	0.826	0.857	0.884	0.885
G4 Reading	0.803	0.820	0.834	0.830	0.871	0.871

183 Table 2 shows that classification consistencies for the predicted scale scores by
 184 XGBoost are higher in all conditions. Mostly, the classification consistencies for the

185 predicted scale scores by Bayesian networks are close to those by XGBoost regression
 186 tree, and much higher than those by linear regression. One exception is for Grade 8
 187 reading test, the classification consistency index for the predicted score by Bayesian
 188 networks at the first performance level cut standard is lower than that by the linear
 189 regression.

190 Prediction Errors

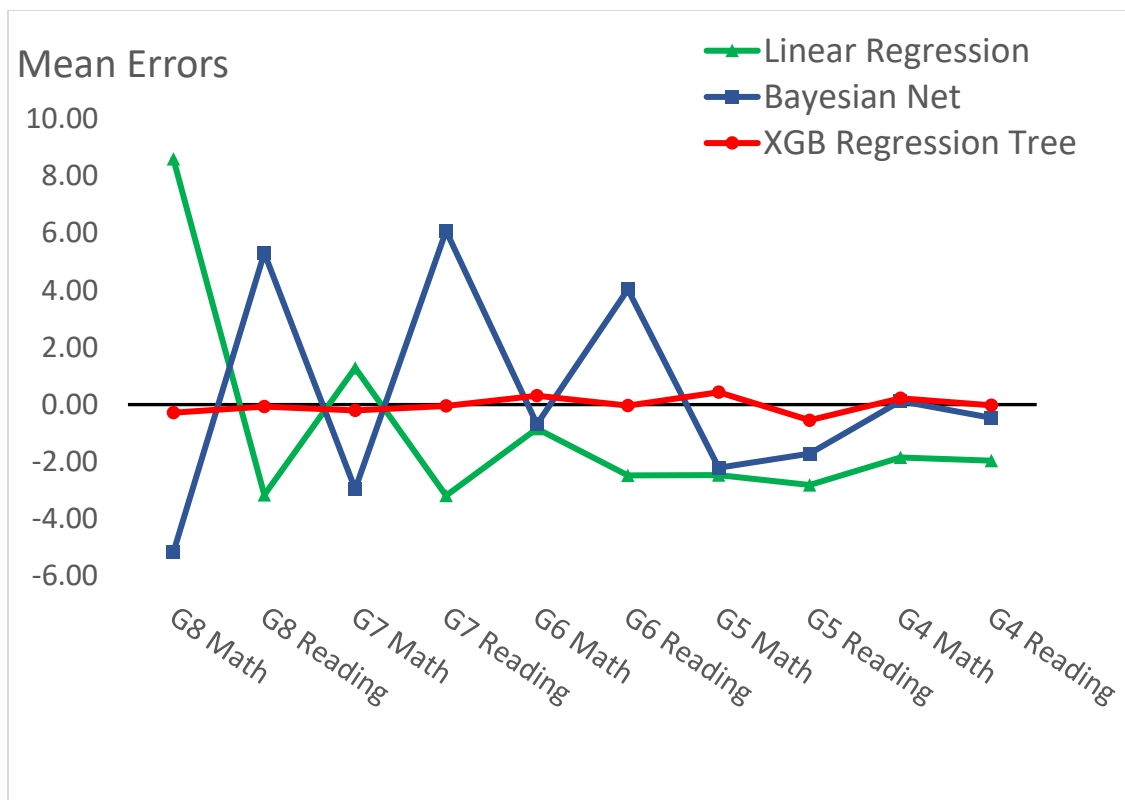
191 The precision of predicted scores by three models was further evaluated using
 192 RMSE. Figure 2 presents RMSE results of three methods.



193

194 *Figure 2* RMSE for all tests by three methods

195 Figure 2 shows that the XGBoost regression tree has the smallest RMSE among
 196 the three methods. Bayesian networks are slightly worse than XGBoost and better than
 197 linear regression for most subjects and grades, except for grade 8 mathematics. In
 198 addition, we also compute the mean errors and find that XGBoost has the most stable
 199 and lowest mean absolute errors across all tests (see Figure 3).



200

201 *Figure 3* Mean absolute errors for all tests by three methods

202 The reason why Bayesian networks don't perform well might be that quite a few
 203 students have missing values for their previous-year scores, and Bayesian networks
 204 would provide bad predictions for these students. On the contrary, XGBoost provides

205 better predictions for students with missing data. In the next section, we conducted
206 some further analysis to test our hypotheses.

207 **The Prediction for Students with Incomplete Data**

208 Generally speaking, students with incomplete inputs have less accurate
209 predicted scores than the students who have complete input variables. Among the three
210 methods, XGBoost regression trees can handle missing data the best, with the highest
211 efficiency. It is able to train models with incomplete datasets and make predictions for
212 incomplete data; The trained model remains stable with or without missing values.
213 Table 3 shows the RMSE for complete and incomplete test datasets respectively, when
214 the XGBoost model was trained with both complete and incomplete data.

215 As a comparison, incomplete data needs to be attended more carefully in
216 Bayesian networks modeling. First, as mentioned above, all functions in 'bnlearn'
217 requires complete data, thus only students with complete data are included in the
218 training data set; Second, variables with only one constant value are removed from the
219 inputs, otherwise parameters will contain zeros and predictions cannot be generated;
220 Third, for students with incomplete data in the test dataset, imputation needs to be
221 carried out for all students to get a prediction; Fourth, when the number of input
222 variables is large (e.g., 117 input variables for Grade 6), the structure learning process
223 becomes extremely computationally demanding. This was one of the reasons why the

224 net structure was fixed in our study, which might not be the best model for imputation
 225 and prediction. Nonetheless, as shown in Table 4, with all the above issues considered,
 226 Bayesian networks can provide adequate predicted scale scores. The model is also very
 227 stable with incomplete data. The existence of incomplete data doesn't exert an influence
 228 on the prediction of students with complete data.

229 Table 3

230 *RMSE for students with complete or incomplete data using XGBoost*

Target Field	N_Train	Complete		Incomplete	
		N_Test	RMSE	N_Test	RMSE
Grade 8 Math	259282	42506	78.4	22315	112.3
Grade 8 Reading	304416	57770	66.7	18335	84.3
Grade 7 Math	263172	52388	67.7	13405	101.0
Grade 7 Reading	290297	58034	67.1	14541	92.9
Grade 6 Math	279032	59056	69.4	10702	93.3
Grade 6 Reading	286567	59875	67.3	11767	98.0
Grade 5 Math	287105	63300	80.2	8477	110.8
Grade 5 Reading	288978	63815	70.6	8430	105.9
Grade 4 Math	287388	67590	82.9	4258	128.5
Grade 4 Reading	287653	67401	78.4	4513	132.5

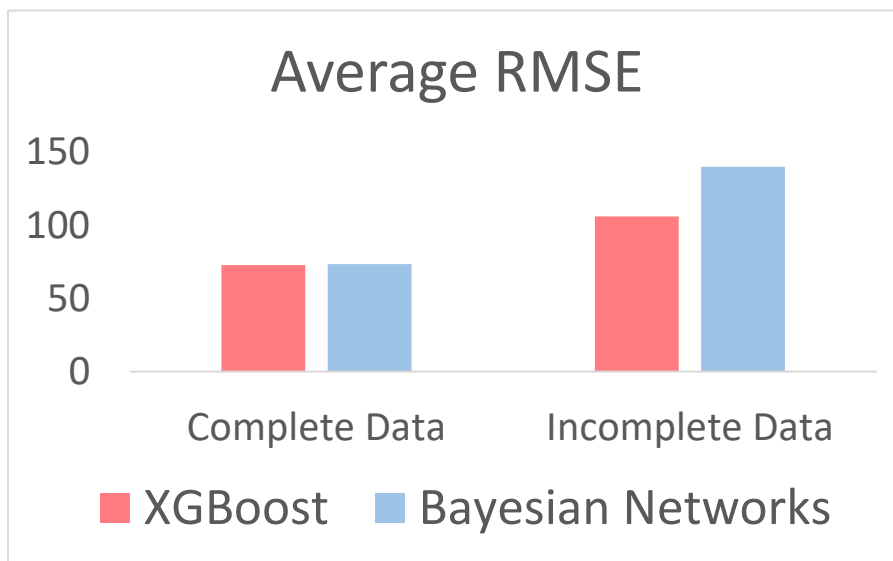
231

232 Table 4

233 *RMSE for students with complete or incomplete data using Bayesian Networks*

Target Field	N_Train	Complete		Incomplete	
		N_Test	RMSE	N_Test	RMSE

Grade 8 Math	136283	42722	79.2	22099	164.6
Grade 8 Reading	185521	58041	66.7	18064	121.7
Grade 7 Math	168379	52624	68.2	13169	118.8
Grade 7 Reading	187311	58334	67.4	14241	143.5
Grade 6 Math	189761	59329	70.9	10429	129.6
Grade 6 Reading	192992	60187	67.6	11455	154.6
Grade 5 Math	204358	63650	80.6	8127	122.9
Grade 5 Reading	205342	64130	70.9	8115	129.1
Grade 4 Math	217450	67965	83.9	3883	153.4
Grade 4 Reading	217133	67816	79.5	4098	160.0



234

235 *Figure 4* The average RMSE across grades and subjects

236 As shown in Figure 4, students with incomplete inputs have less accurately

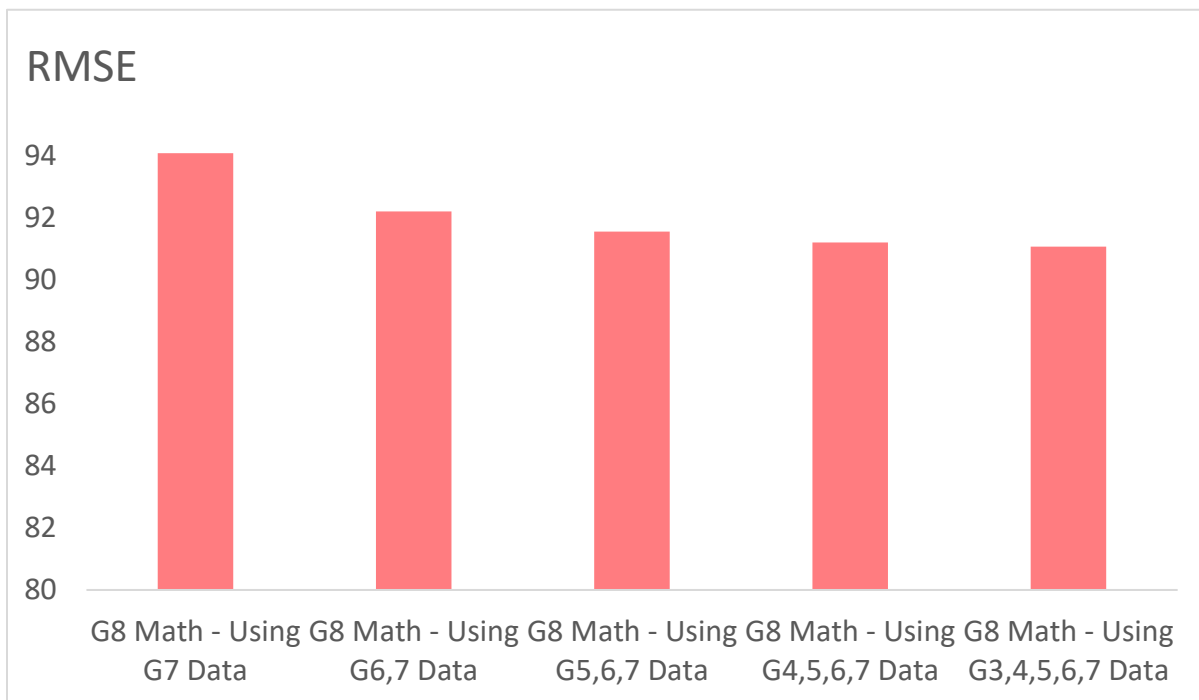
237 predicted scores than the students who have complete input variables. XGBoost handles

238 missing data innately. Specifically, it is able to train models with incomplete datasets;

239 the trained model remains stable with or without missing values. As a comparison,
240 Bayesian networks provide less accurate predicted scale scores for students with
241 incomplete data, even though the missing values were attended more carefully.
242 Nonetheless, the existence of incomplete data doesn't exert an influence on the
243 prediction of students with complete data for both methods.

244 **How Many Previous Years of Data Are Needed?**

245 The prediction errors of XGBoost regression trees using different number of previous-
246 year scale scores are computed. Figure 5 shows that when the number of previous years
247 increased, the prediction accuracy also increased.



248

249 *Figure 5* Decreasing prediction errors with more previous years of data

250

Conclusion

251 The practical purpose of this research is to investigate the practicality of using a
252 statistical framework like XGBoost to forecast scores for next year's tests. The hope is
253 that forecasted scores could then be acted upon by stakeholders, perhaps to identify
254 areas of weakness or focus on at-risk students. In this study, we only predicted future
255 overall scale scores, but the XGBoost statistical framework should be capable of
256 predicting other more specific outcomes, such as more specific test subjects (known as
257 reporting categories in many states).

258 The results indicate that among the 3 statistical approaches (XGBoost, Bayesian
259 Networks, Linear Regression), XGBoost had the best predictive accuracy. This can be
260 expected given the expressive and robust nature of XGBoost, which has proven itself
261 across many big data predictive tasks. In this study, we tuned the XGBoost algorithm
262 specifically for longitudinal test data and were able to successfully create accurate
263 forecasted results. Operationally, XGBoost is very easy to use, as it handles data with
264 missing and incomplete values inherently. Unlike other big data methods, XGBoost
265 offers good interpretive properties as well, enumerating exactly how the model arrives
266 at its output. On the contrary, Bayesian networks require additional considerations in
267 handling missing data, and provide less accurate predictions for students with
268 incomplete data.

269 There are many possible statistical frameworks that could underly models that
270 forecast future performance, and there are almost certainly many additional
271 refinements we could have made to the Bayesian Networks and linear regression
272 models in this study. Our overarching hypothesis, though, is that methods like XGBoost
273 will be able to provide the most accurate predictions even as the number of explanatory
274 variables expand, as expressive models like XGBoost have shown to be very successful
275 across many big data prediction tasks. The results presented in this study can contribute
276 to a fuller understanding of how modern statistical methods can solve or improve on
277 problems of prediction in large-scale measurement.

278

279

280

Reference

- 281 Benjamin, A.S., Fernandes, H.L., Tomlinson, T., Ramkumar, P., VerSteeg, C., Miller, L.,
282 & Kording, K.P. (2018). Modern machine learning far outperforms GLMs at predicting
283 spikes. *Frontiers in Computational Neuroscience*, 12 (56), 1-13.
- 284 Chen, T. & Guestrin, C. (2016). *XGBoost: A Scalable Tree Boosting System*. Paper presented
285 in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge
286 Discovery and Data Mining, San Francisco.
- 287 Friedman, N. (1997). *Learning Belief Networks in the Presence of Missing Values and Hidden*
288 *Variables*. Paper presented in Proceedings of the 14th International Conference on
289 Machine Learning, San Francisco, pp. 125–133.
- 290 Pearl, J., & Russell, S. (2000). Bayesian networks. [online]. Retrieved from
291 <https://escholarship.org/uc/item/53n4f34m>
- 292 Mislevy R. J., Almond R. G., Yan D., Steinberg L. S. (2000). *Bayes nets in educational*
293 *assessment: Where do the numbers come from?* (Tech. Rep. No. 518). Los Angeles, CA:
294 National Center for Research on Evaluation, Standards, and Student Testing.
- 295 Romero, C. & Ventura, S. (2010). Educational data mining: a review of the state-of-the-
296 art. *IEEE Transactions on Systems Man and Cybernetics Part C (Applications and*
297 *Reviews)*, 40(6):601-618.
- 298 Scanagatta, M., de Campos, C. P., Corani, G., & Zaffalon, M. (2015). *Learning Bayesian*
299 *networks with thousands of variables*. Paper presented at 29th Conference on Neural
300 Information Processing Systems, Montreal, Canada.
- 301 Scutari, M. (2010). Learning Bayesian Networks with the bnlearn R package. *Journal of*
302 *Statistical Software*, 35(3): 1-22.

- 303 Sinharay, S. (2006). Model diagnostics for Bayesian networks. *Journal of Educational and*
304 *Behavioral Statistics, 31*, 1-33.
- 305 Tsamardinos, I., Laura E. B., & Constantin F. A. (2006). The max-min hill-climbing
306 Bayesian network structure learning algorithm. *Machine Learning, 65(1)*: 31-78.