



Girls of Steel Robotics

Carnegie Mellon University
The Robotics Institute

CLUTCH MATCHES ARE IN THE MIDDLE

OPTIMIZING OPR WITH WEIGHTED LEAST SQUARES

Gabriel Krotkov

Mentor

Anuva Ghosalkar

Strategy & Scouting Lead

Sienna Li

Data Science Lead

Anushka Prabhu

Data Science Team

Lily Tang

Data Science Team

Audrey Zheng

Data Science Team

Aashi Bhatt

Data Science Team

ABSTRACT

1 In this paper, we present an improvement to Offensive Power Rating (OPR), a popular linear regression
2 model for assessing team performance at a given event. One key assumption of linear regression
3 is the independence of the errors, but in the *FIRST*[®] Robotics Competition (FRC) context, this
4 assumption is not exactly true. Using data from all district events between 2009 and 2024, we model
5 the unweighted errors as a function of tournament progression and generate weightings to improve
6 the regression fit through Weighted Least Squares (WLS). The best weightings show that the most
7 representative matches for a team's overall performance are midway through the tournament. That is,
8 the **real** clutch matches are in the middle.

9 1 Introduction

10 The FIRST Robotics Competition is an international high school robotics competition. Teams build a robot to play a
11 new game released by FIRST each year. In the 2013 game, for example, robots were tasked with shooting frisbees and
12 climbing a seven-foot-tall pyramid. After six weeks of building, teams take their robots to competitions, where they
13 are ranked through a series of qualification matches. Qualification matches consist of six randomly selected robots
14 divided into two opposing alliances of three teams each. To strategize effectively and play to each robot's strengths,
15 teams frequently need to make convenient, flexible, and accurate estimations of their partners' overall performance
16 without relying on large amounts of scouted, robot-level data. OPR uses linear regression to estimate how many points
17 an individual robot contributes to an alliance in any given match. We will refer to a match as a qualification match
18 competed by a particular alliance.

19 1.1 Background

20 Since Karthik Kanagasabapathy and Ian Mackenzie of FRC team 1114 developed "calculated contribution" in 2004,
21 teams have been using linear regression to measure a team's scoring ability. Calculated contribution, generally
22 considered the first application of linear regression in FRC, improves on average score by accounting for alliance
23 partners. The first public description (Weingart,) of similar linear algebra came in 2006 from Scott Weingart of FRC
24 team 293 on a Chief Delphi post which coined the term "Offensive Power Rating". Weingart's terminology for linear
25 regression in FRC rose to popularity, commonly abbreviated as OPR. However, the usual formulation for the regression
26 design matrix uses teams on the columns and matches on the rows (as described by Karthik (Kanagasabapathy,)),
27 where Weingart's formulation used teams on both the rows and the columns. In 2017, Eugene Fang detailed the math
28 behind OPR in a blog post for TheBlueAlliance (Fang,) (TBA), which has become the standard definition of linear
29 regression for FRC.



30 1.2 OPR as Linear Regression

31 By definition, OPR is a multiple linear regression. To see this, reference Section 5.2 (pg 130-133) of Sheather
32 (Sheather,), which shows the calculation behind multiple linear regression. Looking at the solved equation for β ,
33 $\hat{\beta} = (X'X)^{-1}X'Y$, we see that the TBA blog solved equation to find OPR, $x = M^{-1}s$, follows a very similar structure.
34 With β and x as the OPR value, M and $(X'X)$ as the vector of alliance lineups, and s and $X'Y$ as the vector of alliance
35 scores, Sheather's and TBA's explanations are identical. See (Krotkov,) for a demonstration of the equivalence on data.

36 1.3 Motivation

37 Other metrics have so far proven more effective than OPR for match prediction (Statbotics,), in large part because they
38 incorporate historical data. For a simplified example: knowing that FRC team 254 has won 5 world championships
39 is *very* relevant to making a good prediction about their performance in a given match. However, regression does
40 not consider data from previous seasons or previous events. OPR will not see a difference between a multiple world
41 champion and playoff bubble team unless it is given match data from the relevant event. This highlights the primary
42 value of regression methods in FRC: summarizing team performance at a given event. This isn't to say we should
43 disregard match prediction! Match prediction is extremely useful as an empirical test for our methods; but if the goal is
44 accurate match prediction, regression is not the best tool available. To improve the quality of our estimate, we can focus
45 on the model's assumptions and their validity in an FRC context.

46 One of the key assumptions of linear regression is that the errors ε_i are independent and identically distributed with
47 constant variance. That is, the prediction error for each match does not depend on any other matches, and that the spread
48 of the errors does not change over the course of a tournament. In FRC, this is mostly true: matches are well isolated
49 from each other, and the challenges that teams face over the course of the event do not change on average. However,
50 teams do gain experience over the course of the event, make adjustments, and change strategies. Anecdotally, teams
51 "settle in" to their most representative performance after their first few matches, with earlier matches contributing less to
52 overall performance. This dynamic likely influences the distribution of the errors so that it is not constant. Accounting
53 for this non-constant error could improve the linear approximation of team quality.

54 1.4 Weighted Least Squares

55 WLS is a statistical method used to improve regressions by modeling nonconstant residual variance. This method
56 requires the user to know the variance structure of the response in order to model nonconstant variance - or to have
57 a good approximation of it. Typically for cases that cannot be resolved by a transformation, each row of the design
58 matrix is weighted proportional to the inverse variance of the errors (pg. 96, 97 (James, Witten, Hastie, & Tibshirani,)),
59 giving less weight to the less precise observations. This theoretical quantity is usually best approximated by the inverse
60 variance of the residuals. A residual is defined as the difference between the value predicted by a linear regression
61 model and the observed (true) value. The residual in our context is the difference between the score predicted by OPR
62 and the actual alliance score.



63 WLS is generally implemented by dividing the diagonal of the regression covariance matrix by the weights, which
64 directly downweights the least precise observations. This has computational advantages because it deals directly with
65 the regression formulation. However, WLS can be equivalently implemented (Krotkov,) by "row-replication", which
66 duplicates rows an integer number of times to represent the additional weight placed on that row. This equivalent
67 formulation highlights the intuition behind weighted least squares: putting additional importance on each row of the
68 design matrix proportional to the size of the weight. Row-replication also has a flexibility advantage. Applying the
69 weights to the covariance matrix requires a regression setting, while row-replication can be applied without the context
70 of regression.

71 2 Methods

72 2.1 Data

73 Our analysis used qualification match data (alliance lineups and final total scores) from every district event that occurred
74 between 2009 and 2024. Districts have exactly twelve qualification matches for each team, which makes comparison
75 between events more consistent. To ensure that scores are comparable across different years, we standardized the scores
76 of each event.

77 2.2 Weight Estimation

78 WLS allows us to remove the assumption of residual independence, if we can appropriately weight the rows of the
79 design matrix. This means that we need a principled way to find weights that describe the distribution of the residuals.
80 We computed descriptive weights in two ways: residual variance binning and linear weight smoothing. While each way
81 provided a different set of weights, the second (linear weight smoothing), extends the first (residual variance binning).

82 2.2.1 Residual Variance Binning

83 Optimal weights for WLS are proportional to the inverse variance (see (James et al.,) pg. 97) of the error for that
84 data point. To approximate this, we calculate the variance of the residuals of the unweighted linear model in six
85 sequential bins. To bin the residuals, we give each match a "match percentile", which is $\frac{match_number}{n_matches}$; intuitively,
86 this is the percentage of progress through the tournament at which the match takes place. Then, we evenly divide the
87 alliance-matches into six bins based only on their match percentile. Since each team plays twelve matches, six bins
88 allows for an average of two matches from each team to be in each corresponding bin. This balances the granularity to
89 avoid hyper-analyzing the differences or failing to recognize the trend.

90 Taking the variance of the residuals in each bin provides a numerical way to measure the reliability of OPR. To reflect
91 the general trend of the residuals with a unique set of weights, we take the reciprocal of each binned variance.



92 2.2.2 Linear Weight Smoothing

93 Variance binning directly approximates the model's variance; but this could under-smooth and overfit the data. To
94 mitigate this, we also compute linearly smoothed weights, taking the linear regression of the residual variances. Linear
95 weight smoothing significantly reduces the number of weight combinations to test on the OPR model when trying to
96 find the optimal set of weights, providing a more efficient process than residual variance binning in most cases.

97 Linear weight smoothing continues from where variance binning leaves off. By graphing the residual variances against
98 event progression percentile, we can then fit an appropriate function onto the coordinate points. With the residual
99 variances "smoothed" with linear approximation, we then obtain the six variances (one for each bin) projected by
100 residual variance graph. The reciprocals of the projected residual variances join to form a set of weights unlike the
101 weights provided by variance binning.

102 2.3 Weight Evaluation

103 To find the best weights, we are interested in risk, a model's error on the test data, rather than its loss, a model's error on
104 the training data. We utilized two methods to evaluate the weighted OPR model, both of which make judgments based
105 on the risk produced by both the weighted and unweighted OPR models. The key difference between the two methods
106 is *how* the risks are calculated.

107 For both methods, to contextualize the performance of the weighted model against the unweighted model, we take the
108 ratio $\frac{R_{unweighted}}{R_{weighted}}$, where R is some estimate of a model's risk. This represents how much the weighted model improves
109 on the unweighted model; since a higher risk is worse (higher error), we can interpret this ratio as a proportion.

110 2.3.1 Test Mean Squared Error

111 Mean squared error (MSE) is a measure that can adequately characterize the accuracy of a model. We calculate MSE by
112 taking the average of the squared residual variances.

113 MSE only measures the error of a model applied on a fixed and complete dataset. Consequently, it becomes vulnerable
114 to overfitting and does not evaluate models based on their ability to make predictions for unseen data. We are interested
115 in using methods that do a better job of estimating our model's *prediction risk*, rather than just its in-sample MSE.
116 Test-set cross validation is a method that avoids these issues by splitting data into a training and testing set; we train the
117 model only on the training data, reserving the testing data to evaluate the model.

118 Combining both of these methods leads to test MSE, which is the MSE of only the test set. However, withholding a test
119 set sacrifices a decent proportion of the complete data for training, resulting in high variance.

120 2.3.2 Leave-One-Out Cross Validation

121 Leave-One-Out Cross Validation (LOOCV) is another effective risk estimate from machine learning. Like test MSE,
122 LOOCV splits the data into a test and training set, but instead uses only one data point as the test set. This method



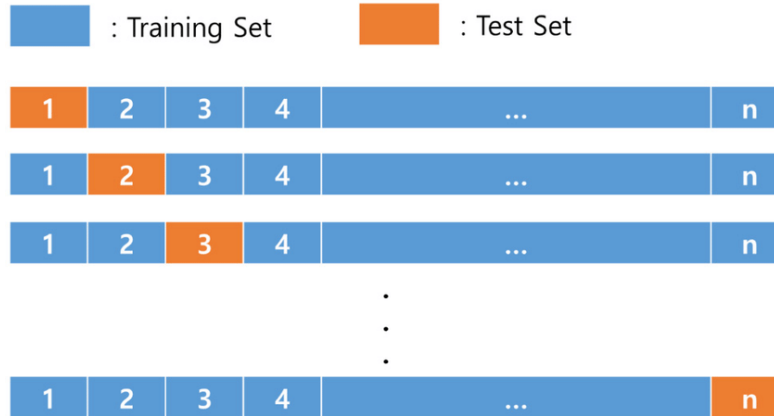


Figure 1: Leave One Out Cross Validation Diagram

123 trains and produces a linear regression model using all but one of the matches in an FRC event (e.g. in a 160-match
 124 event, 159 matches are used for training). It then tests the trained model by predicting the omitted match's score and
 125 recording the residual. LOOCV repeats this process for each match in the event so every match is tested once. We
 126 prefer LOOCV over test MSE because it uses nearly all of the data, therefore maintaining low variance. LOOCV is
 127 detailed visually in 1 and mathematically in pages 200-203 of (James et al.,).

128 2.4 Weight Optimization

129 To pick the number of bins we used hyperparameter tuning with the mean LOOCV as the loss function. However, once
 130 you have b , picking the optimal weight in each slot poses a computationally difficult challenge, in $O(n^b)$. This makes
 131 an exhaustive grid search difficult, so we attempted three solutions to optimize an easier problem and approximate the
 132 optimal weights. First we considered single-bin optimization, a strategy that only changed one from the originally
 133 estimated weights. Then we considered a stepwise simultaneous solution which optimizes each bin in the context of the
 134 originally estimated weights and then picks the independently optimized value for each bin. Finally, we considered
 135 sequential weight fixing, which optimizes in a single bin and fixes that optimized value for the rest of the optimization.

136 3 Results

137 3.1 Exploratory Data Analysis

138 Figure 2 shows the linear relationship between the qualification match percentile and the resulting squared residual.
 139 The line of best fit has a slope of -0.031 ($p < 0.001$), demonstrating that qualification matches that occur later in the
 140 tournament have smaller squared residuals. The negative slope means that on average, the error decreases as match
 141 percentile increases, which confirms a measurable relationship between match percentile and the residual produced.
 142 This supports our assertion that weighting reduces variance by determining how important each match is and assigning



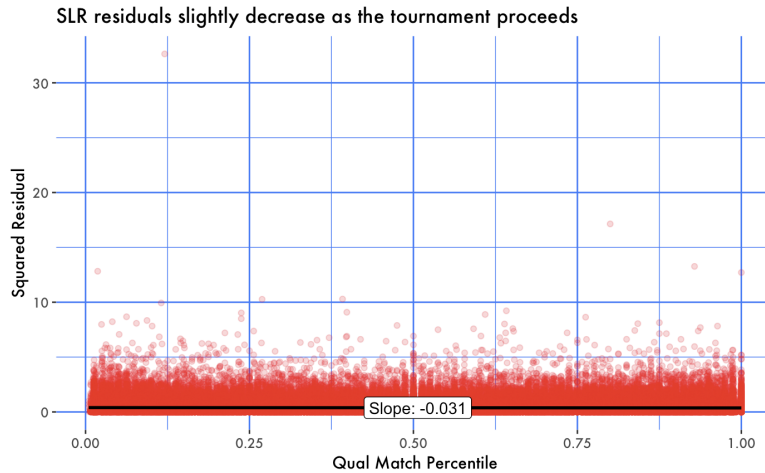


Figure 2: Simple Linear Regression on Raw OPR Residuals by Match Percentile

143 its appropriate relevance in calculations. Figure 2 confirms the advantage of weighting late-tournament matches more
 144 heavily by observing how they have slightly less variance overall, thus reducing error and improving predictions.

145 3.2 Weight Estimation

146 In Figure 3a, the residual variance is plotted for each of the six bins. These values represent the spread of error within
 147 each bin. The smaller this spread, the more consistent the unweighted OPR residuals. Overall, we see a parabolic
 148 trend: the residual variances are very high towards the beginning (bins 1 and 2), low towards the middle (bins 3 and 4),
 149 and increase somewhat again towards the end of the tournament (bins 5 and 6). This means that the unweighted OPR
 150 estimations are the least consistent at the extremes of a tournament, especially at the beginning. Hence, OPR should
 151 be given more weight for matches that fall in bins 3 and 4, and less weight for matches that lie further towards the
 152 beginning or end of the tournament. Linear weight smoothing produces a v-shaped variance graph with the reciprocals
 153 of the estimated weights, as shown in 3b.

154 3.3 Weight Optimization

155 Tuning over the number of bins between 2 and 15 found 12 to be the ideal number of bins, as shown in 4d. This makes
 156 intuitive sense because this data is based entirely on district performances, which each have 12 matches. Dividing
 157 span of matches into 12 bins approximately mimics the progress of the tournament in "rounds" where each robot plays
 158 approximately once.

159 Stepwise optimization performed the best of our optimization strategies, achieving a mean LOOCV of 0.666 as opposed
 160 to single-bin optimization's 0.667 and sequential weight fixing's 0.669. The optimal weights found were the following:

$$w_{stepwise} = (1.90, 2.56, 2.94, 3.13, 3.24, 3.20, 3.17, 3.15, 2.92, 2.79, 2.66, 2.43)$$



Figures: Weighting Estimation & Evaluation

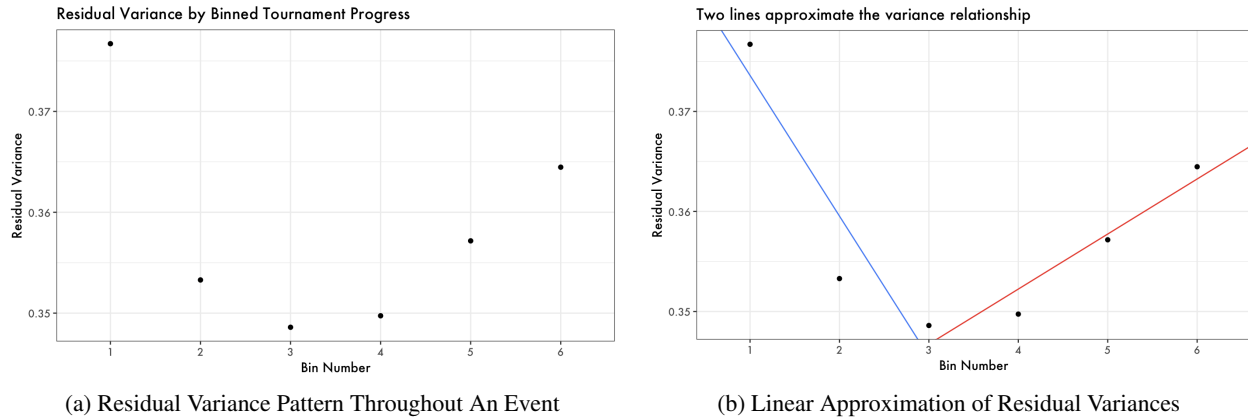


Figure 3: Residual Variance Patterns

161 3.4 Test Mean Squared Error

162 Recall from section 2.3.1 that we calculate Test MSE by producing a model based solely on the training data and
 163 computing that model's MSE from the testing data. In figure 4a we plot the ratios between the unweighted and weighted
 164 event test MSEs. The mean of the distribution is 1.004, indicating that weighted OPR is .4 percent more accurate than
 165 unweighted OPR on average. The 95% confidence interval around the mean value is computed by bootstrap(Lomuscio,
 166).

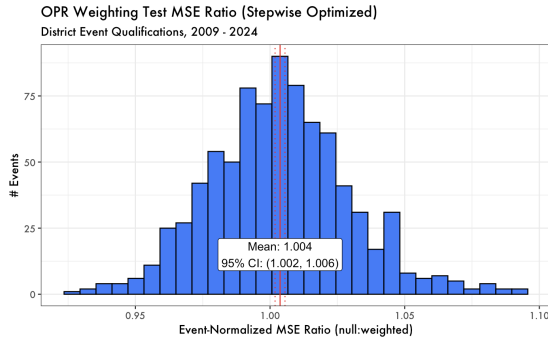
167 3.5 Leave-One-Out Cross Validation

168 An Event-Normalized LOOCV Error Ratio above 1 suggests that weighted OPR has a measurable predictive advantage
 169 over unweighted OPR. Figure 4b shows a mean LOOCV ratio of 1.004, suggesting that on average, weighting does
 170 improve OPR predictions by 0.4% as directly compared to its unweighted counterpart. The 95% confidence interval
 171 around the mean value is computed by bootstrap(Lomuscio,).

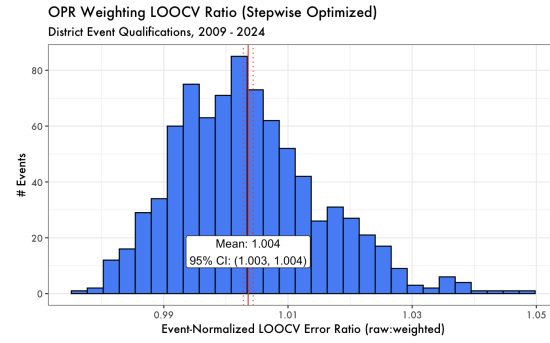
172 Figure 5 shows the LOOCV ratio broken down by year. In nearly all cases, the binned LOOCV ratios are larger than
 173 the linear ratios, confirming that bin weighting outperforms linearized weights. Furthermore, 2010 is the only year
 174 where the binned LOOCV ratios are below 1. Although this implies that raw OPR possesses an advantage here, there
 175 were only eight events run this year - the small sample size makes it difficult to make confidence inferences based on
 176 that year. Stepwise optimization extracts improved performance from games favorable for OPR (with linear, separated
 177 scoring like 2019, 2022, and 2016) at the cost of worse performance when the key assumptions of OPR are invalid, like
 178 in 2014 or 2018.

179 Figure 4c shows the percentage improvement of weighted OPR over raw OPR over time. The size of each dot scales
 180 with the number of events for that year recorded in the data, which is the number of district events played that year. The
 181 blue dotted line is the average of the data. The only year with a negative percent improvement was 2010; however, note

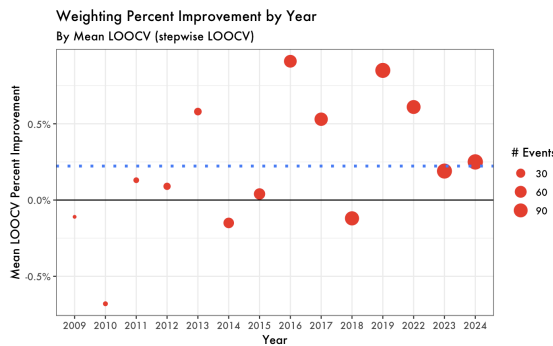




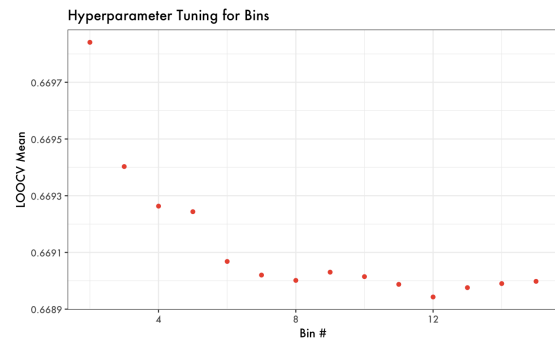
(a) Test MSE Ratio



(b) Leave One Out Cross Validation Ratio Distribution



(c) Weighting Percent Improvement over Raw OPR



(d) Tuning for the number of bins

Figure 4: Weighting Evaluations

182 that before 2014 no year had more than 20 district events. As the number of events grows, we see the improvement
 183 from weighting stabilize on average.

184 **4 Discussion and Conclusion**

185 Independent residuals is a key assumption for linear regression. Under this assumption, linear regression is extremely
 186 efficient. However, we proved that the residuals of the OPR model are in fact not independent. Taking advantage of this
 187 invalid assumption, we can lean into the matches in which OPR is most predictive, leading to an improved model. The
 188 weightings that performed the best followed an asymmetric, roughly parabolic shape, with by far the least weight on
 189 very early matches and the most weight on matches midway through the tournament.

190 Weighting OPR based on match recency showed a consistent, but small improvement over unweighted OPR at district
 191 events between 2009 and 2024. Weighted Least Squares improved our estimation of teams by 0.015 Crescendo points
 192 and impacted teams' OPRs by about 0.2 points. This is not a large change, but it shows that even an unoptimized
 193 weighting can improve on unweighted OPR.

194 We found that the variance of unweighted OPR residuals follows a similar pattern in every year we tested 5. Notably,
 195 weighting improves OPR *independent* of OPR's value as a metric in a given year. For example, OPR did very well in
 196 2022 and poorly in 2017, but the accuracy of OPR increased similarly with weighting for both years 4c.



Mean LOOCV Ratios				
unweighted:weighted ¹				
Year	# Events	Linear	Binned	Stepwise
2009	7	1.0005	1.0007	0.9989
2010	8	0.9987	0.9993	0.9932
2011	10	1.0001	1.0009	1.0013
2012	17	1.0002	1.0008	1.0009
2013	19	1.0011	1.0019	1.0058
2014	44	1.0003	1.0003	0.9985
2015	53	1.0003	1.0007	1.0004
2016	72	1.0021	1.0025	1.0091
2017	81	1.0016	1.0018	1.0053
2018	94	1.0004	1.0004	0.9988
2019	109	1.0021	1.0024	1.0085
2022	90	1.0016	1.0020	1.0061
2023	109	1.0008	1.0010	1.0019
2024	114	1.0009	1.0011	1.0025

¹ Higher values indicate the weighting is performing better

Figure 5: LOOCV Ratios between Unweighted and Weighted OPR

197 4.1 Limitations

198 WLS is only as good as its weights. We identified *good* weightings using stepwise optimization, but did not numerically
 199 optimize to find the *best* weightings over a full grid. For our first weighting model, residual variance binning, the
 200 computational cost to optimize the weighting over b bins is in $O(n^b)$. This would take around four days of computation
 201 over a reasonable grid, which was outside our scope.

202 To make optimization easier, we tried weight linearization (3b), which simplifies our search space. Instead of optimizing
 203 over b bins, we would optimize over two slopes and two intercepts, for $O(n^4)$. However, this simplification trades
 204 accuracy for optimization speed. In both cases, more optimization is required before we find the numerically *best*
 205 weights.

206 4.2 Next Steps

207 The next step to improving our weighting is tuning the slopes and intercepts of the piecewise linear function using cross
 208 validation. With considerable computing resources, it would also be worthwhile to brute-force optimize over the binned
 209 weights.

210 In a broader scope, the primary limitation of linear regression in an FRC context is small sample size. Teams almost
 211 never play more than twelve qualification matches each in a single tournament, and while simple linear regression is an
 212 efficient estimator, it doesn't stabilize until late in the tournament. To make a more useful single-event summary for
 213 FRC, we need to continue to reduce the variance, to make a model that stabilizes faster. However, per the Gauss-Markov
 214 theorem (Taboga,), we know that simple linear regression is the minimum variance unbiased linear estimator. Therefore,



215 to improve on the variance of that estimator, you either need to accept bias, similar to how ELO models incorporate
216 historical information, or adopt a nonlinear model. Both options provide promising avenues to improving on current
217 regression methods in FRC.

218 **4.3 Applications Outside the Regression Context**

219 Any set of weights can be used to create a design matrix with repeated data entries, where the weighting vector would
220 determine how many times a set of matches within a bin should be duplicated. Running simple linear regression on this
221 matrix is equivalent to running WLS on the original match matrix.

222 While WLS would use a set of optimized weights to compute OPR coefficients immediately, row replication would
223 generate a design matrix with duplicated rows before applying simple linear regression to output the *same* coefficients.
224 Using row replication allows us to apply this data in a nonregression context and use models that incorporate higher
225 biases to find more accurate results.



226 **5 Technical Appendix**

227 Code for implementing weighted least squares for OPR can be found in the scoutR repo, at
228 scoutR/markdown/opr_weighting. To create the data file district_qualms_09_24.rda, use this script.

229 Code for hyperparameter tuning can be found here.

230 **6 Contact**

231 To reach out to the Girls of Steel Data Science team, email Girls of Steel.

232 For questions or for help replicating our work, email Gabriel Krotkov.

233 **7 Acknowledgments**

234 We would like to express gratitude to the reviewers of our paper who gave us constructive feedback during our writing
235 process, Andy Miller-Peterson, Xiaohan Liu, and Karthik Kangasabapathy. We would also like to thank the Girls of
236 Steel mentors who supported the Girls of Steel data science team meetings over the past 5 months, namely Abby Shrack,
237 Joe Jackson, and Andy Miller-Peterson.



238 **References**

- 239 Fang, E. (2017). The math behind opr - an introduction. *TheBlueAlliance Blog*.
- 240 James, G., Witten, D., Hastie, T., Tibshirani, R. (2021). *An introduction to statistical learning*.
241 Springer.
- 242 Kanagasabapathy, K. (2020). Effective first strategies. In *Rsn spring conferences presented by wpi*.
- 243 Krotkov, G. (2024a). Opr as linear regression. <https://tinyurl.com/opraslinreg>.
- 244 Krotkov, G. (2024b). Weighted least squares and row replication. <https://tinyurl.com/27bc6565>.
- 245 Lomuscio, S. (2023). Bootstrap estimates of confidence intervals. *University of Virginia Library*.
- 246 Sheather, S. J. (2009). *A modern approach to regression with r*. Springer Texts in Statistics.
- 247 Statbotics. (2023). Evaluating frc rating models. <https://www.statbotics.io/blog/models>.
- 248 Taboga, M. (2021). Gauss markov theorem. *Lectures on probability theory and mathematical*
249 *statistics*. Kindle Direct Publishing. Online appendix. [https://www.statlect.com/fundamentals-](https://www.statlect.com/fundamentals-of-statistics/Gauss-Markov-theorem)
250 *of-statistics/Gauss-Markov-theorem*.
- 251 Weingart, S. (2006). Let's do a little linear algebra. [https://www.chiefdelphi.com/t/offense-defense-](https://www.chiefdelphi.com/t/offense-defense-rankings-for-1043-teams/71490/19)
252 *rankings-for-1043-teams/71490/19*.

