

A Bayesian Approach to March Madness Prediction

Mohamad Daniel Bairakdar Jacob Cohen Hoang Luong Bryce Shellow

Abstract

In this work, we train a Bayesian logistic regression model on past college basketball season statistics to predict the outcome of the 2022 season March Madness games. We base our predictions on multiple factors, including a variety of regular season box score statistics, ordinal team rankings from multiple ranking systems, and March Madness seedings. We investigate the effects of including box score statistics from less recent seasons at the expense of ranking data, and explore the sensitivity of our model to prior specification for the different ordinal ranking systems. Our best model achieves an average cross validation accuracy of 66.4%, and an accuracy of 58.0% when tested on the 2022 season March Madness games.

1 Introduction

Sports analytics has been a vastly growing field, especially over the past few years. The analytics began in baseball with the Moneyball A's in the early 2000s, and has now expanded to almost all other sports. These include other team sports such as basketball, football, and soccer; individual sports like tennis and golf; and combat sports such as boxing and the UFC. We decided to attempt to predict the outcome of individual March Madness games. March Madness - a 68-team tournament that takes place at the end of the Division I basketball season - is a statistical goldmine with more possible outcomes than one can count. For every season, millions of people across the country make their brackets, but few are able to accurately predict the outcomes. In fact, there has never been a perfect bracket in the history of the tournament. The longest someone has ever gone with a perfect bracket is only 49 of the 63 games. This is in large part due to the fact that there are 9.2 quintillion possibilities for it. While we did not attempt to predict the bracket, we took a step in the direction of predicting it by attempting to predict the winner of any possible matchup between two tournament teams. In theory, once a strong model of this kind is built, it could be stacked sequentially to predict the bracket. In order to analyze each tournament matchup, we used regular season box score statistics, ordinal team rankings, and March Madness seedings to see if they were indicative of March Madness outcomes. We

fitted a Bayesian logistic regression using an MCMC computing algorithm and used cross-validation as a diagnostic tool to pick the best model. We then tested our model by predicting the outcomes of the 2022 season March Madness games. Ours is not the first work to attempt to predict the outcome of college basketball games. Brown (2019) builds a logistic regression for this purpose, but does not use a Bayesian approach. Nelson (2012) employs an approach similar to ours, but only uses box score statistics to predict the outcome, does not investigate the effect of adding or removing data from less recent seasons - data which could impact the prediction since the nature of the game has changed over the years - and does not investigate model sensitivity to prior specification.

2 Data

2.1 Data Source

We obtained our data from the Kaggle dataset “March Machine Learning Mania 2022 - Men’s”¹, which contains single game regular season and tournament box score data dating back to the 1984-1985 college basketball season. Detailed box score data is only available starting from the 2002-2003 season, so we discard the data from earlier seasons. Due to the 2020 Covid-19 pandemic, the National Collegiate Athletic Association (NCAA) canceled the tournament for that year, and thus we did not include 2019-2020 regular season data as the corresponding tournament data was missing.

2.2 Covariates: Regular season data

The regular season data consisted of individual game logs containing variables indicating the year of the season, the day the game occurred in relation to the start of the season, the location of the game, number of overtimes that occurred in the game (if any), the team ID for the winning and losing teams, and their respective box score statistics. These statistics include the number of points scored, the number of field goals (FG) made and attempted, 3 point FG made and attempted, free throws (FT) made and attempted, offensive and defensive rebounds, assists made, turnovers committed, blocks and steals made, and personal fouls committed. The same data was available for the tournament. The Kaggle dataset also contained ordinal rankings from multiple basketball metric sites. For each team, this data included the season, ranking system name, ordinal rank, and day of the ranking since the start of the season - rankings are provided on a weekly basis. Finally, the dataset included the tournament seeding for each team. Figures 1a and 1b show samples of raw box score statistics and ranking data, respectively.

To prepare the data to be input into our model, for each season, we discard data for teams that did not participate in March Madness for that season. Next, we combined regular season team box

¹<https://www.kaggle.com/c/mens-march-mania-2022/data>

```

Season DayNum WTeamID WScore LTeamID LScore WLoc NumOT WFGM WFGA WFGM3
<dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <chr> <dbl> <dbl> <dbl>
1 2003 10 1104 68 1328 62 N 0 27 58 3
2 2003 10 1272 70 1393 63 N 0 26 62 8
3 2003 11 1266 73 1437 61 N 0 24 58 8
4 2003 11 1296 56 1457 50 N 0 18 38 3
5 2003 11 1400 77 1208 71 N 0 30 61 6
6 2003 11 1458 81 1186 55 H 0 26 57 6
# ... with 23 more variables: WFGA3 <dbl>, WFTM <dbl>, WFTA <dbl>, WOR <dbl>,
# WDR <dbl>, WAST <dbl>, WOT <dbl>, WSTL <dbl>, WBLK <dbl>, WPF <dbl>,
# LFGM <dbl>, LFGA <dbl>, LFGM3 <dbl>, LFGA3 <dbl>, LFTM <dbl>, LFTA <dbl>,
# LOR <dbl>, LDR <dbl>, LAsT <dbl>, LTO <dbl>, LStL <dbl>, LB1k <dbl>,
# LPF <dbl>

```

(a) Box score statistics

```

Season RankingDayNum SystemName TeamID OrdinalRank
<dbl> <dbl> <chr> <dbl> <dbl>
1 2003 35 SEL 1102 159
2 2003 35 SEL 1103 229
3 2003 35 SEL 1104 12
4 2003 35 SEL 1105 314
5 2003 35 SEL 1106 260
6 2003 35 SEL 1107 249

```

(b) Ordinal team rankings

Figure 1: Raw data samples

score winning and losing statistics by grouping by team ID and season to create “for” and “against” season averages. For certain variables like field goals attempted and made, we compute average field goal percentage, rather than an average for each of field goals attempted and made separately. We then created variables prefixed with “current” or “other.” “current” indicates the “for” values for the team - in other words, the numbers that the team produced throughout the season - while “other” indicates “against” values, which are the averages and percentages produced by the team’s opponents throughout the season. To obtain individual training instances for our model - where each instance consists of a March Madness matchup between two teams - for each variable we create, we compute a difference between the values for those variables for each team participating in a given game. For example, if Michigan is the first team listed in a Michigan-Kansas 2016 tournament match-up, and Michigan had an average FG% of 0.53 while Kansas had a FG% of 0.43 over the 2015-2016 regular season, then the covariate representing the difference in “current” FG% would be 0.10. We note that one of our variables represents difference in game winning percentage between the two teams pitted against each other; we did not include the difference in the losing percentage, as that would just be the opposite of the difference in the winning percentage.

As explained later, we were interested in evaluating model performance for models built using data spanning different timeframes; a model built using data from the 2016-2021 seasons versus a model built using the data from the 2003-2021 seasons. The ranking information we obtained depended on the number of seasons we used to build our model. Within a given timeframe of interest, there were rankings from certain ranking systems that were available for one season but not another. In those cases, we excluded those ranking systems. Other ranking systems had rankings for some teams but not others. We also excluded those ranking systems. An alternative approach would be to impute those values. Once we ended up with a list of ranking systems that met the above criteria, we then decided based on our prior knowledge of the quality of those rankings which ones to include and which to exclude. The ones we include in our models are: B Wilson Empirical (BWE), Massey (MAS), Pomeroy (POM), and Stat Fox (SFX). For each ranking system, we create two covariates: one that represents the ordinal ranking of the team over the entire regular season, and one that represents

the ordinal ranking right before the start of the tournament. This is to account for the fact that certain teams may be very good teams on average, but may be “out of form” during certain periods of time for many different reasons - most notably because of injuries of key players. In table 1, average rankings over the entire season correspond to variables with the substring “Mean.” Lastly, the variable “SeedAdd” was created as an interaction term to account for the difference in competition between teams; there is a larger performance gap between a 9 seed versus a 16 seed than a 1 seed versus an 8 seed. In this case, both scenarios would result in a “SeedDiff” value of 7, whereas we could tell from the lower “SeedAdd” that the 1 and 8 seeds match more closely in skill than the other match-up.

All variables we created are shown in table 1, with their descriptive statistics for the 2016-2021 timeframe. See the Appendix, Section 8.2 for the code used to create the variables. The variables were normalized by subtracting their respective means and standard deviations.

Variable	Mean	Standard Deviation
winPctDiff	0.004052	0.1519
currScoreAvgDiff	0.7782	7.3043
otherScoreAvgDiff	0.3279	6.4514
currFGPctDiff	0.003390	0.0348
otherFGPctDiff	-0.0008	0.02852
currFG3PctDiff	0.003107	0.03848
otherFG3PctDiff	0.0008	0.0306
currFTPctDiff	0.0015	0.0519
otherFTPctDiff	-0.0036	0.0335
currRebAvgDiff	0.3975	3.9076
otherRebAvgDiff	-0.0555	3.3806
currAstAvgDiff	0.1906	2.7380
otherAstAvgDiff	0.08557	2.1618
currTOAvgDiff	0.04426	1.9082
otherTOAvgDiff	-0.1649	2.4988
currDefAvgDiff	-0.0884	2.3694
otherDefAvgDiff	0.0177	1.6250
SeedDiff	-0.5689	7.4985
BWEDiff	-6.4641	74.1355
MASDiff	-6.4671	71.4674
POMDiff	-6.9910	74.8927
SFXDiff	-6.0479	70.4773
POMMeanDiff	-7.7559	81.9615
BWEMeanDiff	-7.2479	79.8221
MASMeanDiff	-7.8061	80.4986
SFXMeanDiff	-6.4981	74.7913
SeedAdd	14.4132	5.8783

Table 1: Descriptive statistics for all variables we created for the 2016-2021 seasons

2.3 Outcome: March Madness data

From the tournament data, we only obtained the IDs for the winning and losing teams for each game. The tournament data was used to create our response variable, so the game statistics are not needed as we are only attempting to predict the outcome based on the regular season values. In order to obtain

a balanced dataset, within each March madness tournament, we randomly shuffle the games and split the resulting data in half so that the first half consists of games considered as “Wins,” while the other half consists of games considered as “Losses”.

3 Model

3.1 Modeling Approach

We model the problem of predicting the outcome of a March Madness game as a binary task using Bayesian Logistic Regression. When given two teams, Team i and Team j , our model answers the question: will Team i beat Team j when they are pitted against each other? If the answer is yes, the model predicts 1, and if not, it predicts 0. Note that the order in which the teams are input matters. That is, if the same teams i and j are input into the model but in the reverse order, then the model should predict the opposite of what it predicts the first time. Formally, we denote by y_{ij} the indicator variable that represents a win by team i when pitted against team j . If $y_{ij} = 1$, this means team i wins against team j , while 0 means team i loses. We model y_{ij} as follows:

$$\mu_{ij} = \text{logistic}\left(\alpha + \sum_{n=1}^N \beta_n x_{n_{ij}} + \beta_{ss} z_{ss_{ij}}\right) \quad (1)$$

$$y_{ij} \sim \text{Bernoulli}(\mu_{ij}) \quad (2)$$

Here, *logistic* is the logistic function $f(x) = \frac{1}{1+e^{-x}}$, N is equal to the number of covariates, α is the intercept term, and $x_{n_{ij}}$ is the difference between the average of variable x_n for team i and the average of variable x_n for team j . Finally, as discussed at the end of Section 2.2, we add an interaction term $z_{ss_{ij}}$, which is the sum of the seeds for teams i and j .

To fit our models, we use an MCMC algorithm using the R package Stan. Specifically, we use 4 chains that run for a total of 10000 iterations each, with 5000 warmup iterations.

3.2 Model specification

We build three models using the modeling approach just described. The essential characteristics of each model are shown in table 2. For Model 1, we investigate the effect of adding more historical regular season box score statistics at the expense of ranking data, as explained in Section 2.2. Given that we have limited prior information on which covariates are important for the prediction, we impose a very weakly informative prior on all regression coefficients and the intercept. For Model 2, we use only data from the 2016-2021 seasons. As explained previously, since the nature of the game has changed drastically over time, we hypothesized that this may have a significant effect on the model built. Additionally, using less historical data has the advantage of allowing us to use ordinal ranking

data, besides team seeds of course. Again, for Model 2, we impose a very weakly informative prior on all regression coefficients and intercept. Finally, Model 3 is similar to Model 2, except we impose a $\text{Normal}(0, 1)$ prior on the coefficients corresponding to SFX and MAS rankings to induce shrinkage, since we believe the other two ranking systems are of better quality. This will allow us to investigate the sensitivity of our model to prior specification for ordinal rankings.

Model 1	Model 2	Model 3
2003-2021 Seasons	2016-2021 Seasons	Similar to Model 2
No Ordinal Rankings	Ordinal Rankings: BWE, POM, MAS, SFX	Generic weakly informative prior on SFX and MAS: $\text{Normal}(0, 1)$
Very weakly informative prior on all regression coefficients and intercept: $\text{Normal}(0, 10)$	Very weakly informative prior on all regression coefficients and intercept: $\text{Normal}(0, 10)$	Very weakly informative prior on remaining regression coefficients and intercept: $\text{Normal}(0, 10)$

Table 2: Characteristics of the three different models we built

4 Diagnostics

To determine which of our three models would serve as the best predictor, we performed cross-validation. For Model 1, we divide the dataset into 6 folds, where each fold is composed of data from three consecutive seasons, and perform 6-fold cross-validation. For Models 2 and 3, we perform 5-fold cross-validation, where each fold consists of data from a single season. The results are listed in table 3.

Model	Cross validation accuracy (%)
1	63.5
2	66.4
3	66.0

Table 3: K-fold cross validation accuracy for each model

To compute the accuracy, we repeatedly draw samples of each of the model’s parameters from their posterior distributions, compute μ_{ij} from Equation 1 for each game in the validation set, and then sample y_{ij} as described in Equation 2 for each game. Accuracy is then computed by counting the number of correct predictions made by the model over the total number of predictions made. This process is done repeatedly as different parameter values are sampled from their joint posterior distribution, so that in the end we obtain an MCMC accuracy estimate, which is just the mean of the posterior accuracies generated in the MCMC process. In practice, this is done using the Stan function “gqs”, which will compute the posterior distributions of the quantities specified in the “generated_quantities” block of a Stan model (see the “generated_quantities” block of our Stan model in the Appendix, Section 8.3

for more details).

We observe that Model 2 outperforms the other two models, with a 5-fold cross validation accuracy of 66.4%. We thus proceed to train a new model which uses the specification of model 2, trained on the whole 2016-2021 dataset. We check this model for convergence by examining traceplots for all the model's parameters. This is shown in figure 2. We observe that all parameters converge.

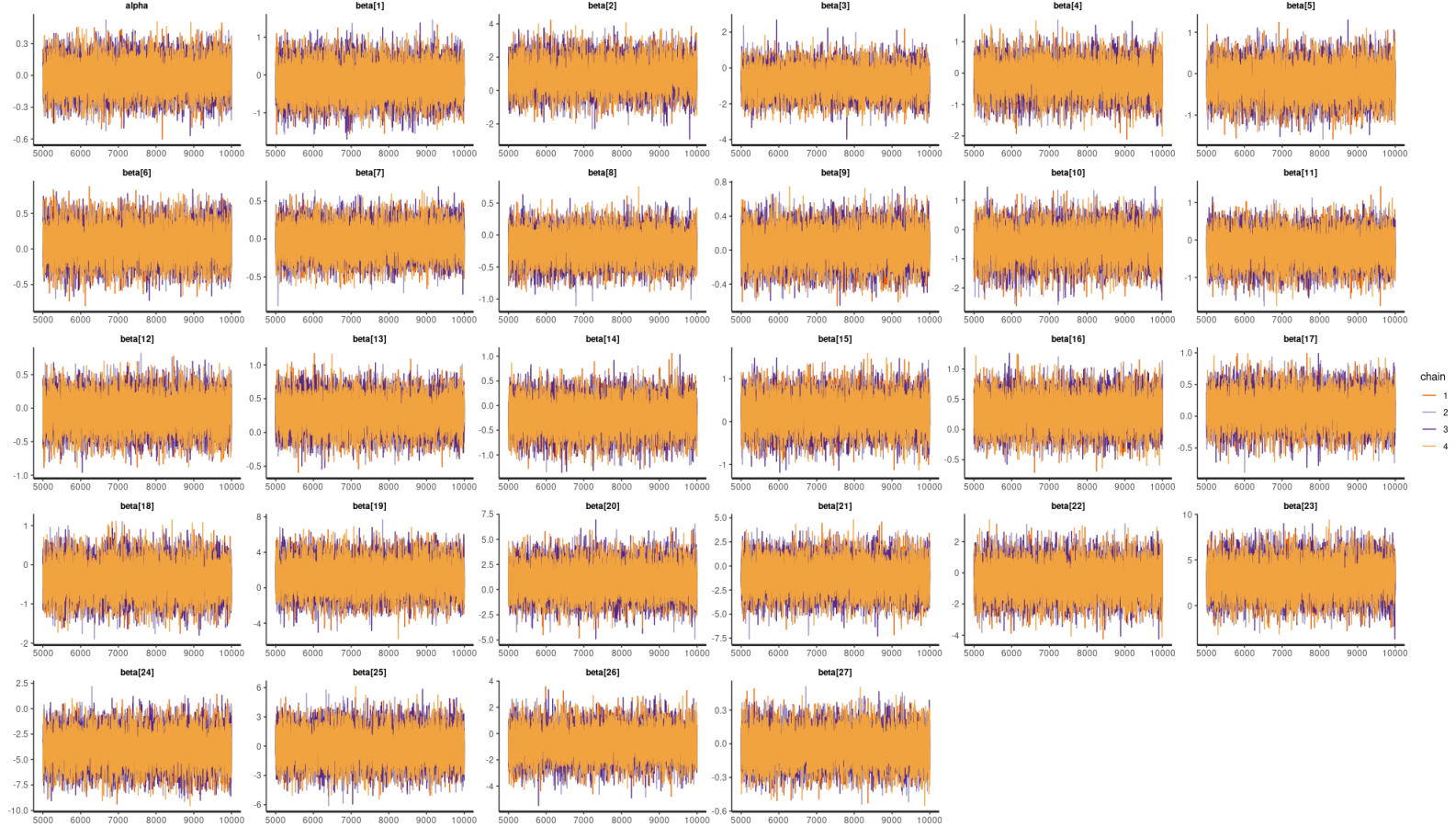


Figure 2: Traceplots for the best model's parameters. Barring the intercept term alpha, the order in which the parameters are listed corresponds to the order in which the variables are listed in table 1

Having checked convergence, we then use our model to make predictions for the games of the 2022 season. We scale all variables for the test set using the training set variables' mean and standard deviations. We obtain an accuracy of 58%. The accuracy for the test set is computed in the same way as it was computed in the case of cross-validation, described above.

5 Model interpretation

We next investigate the 80% and 90% credible intervals for model parameters. These are shown in figure 3. We note that except for beta[23] and beta[24], which correspond to the covariates POMMeanDiff and BWEMeanDiff, the 80% (and hence the 90%) credible intervals of all other parameters contain 0.

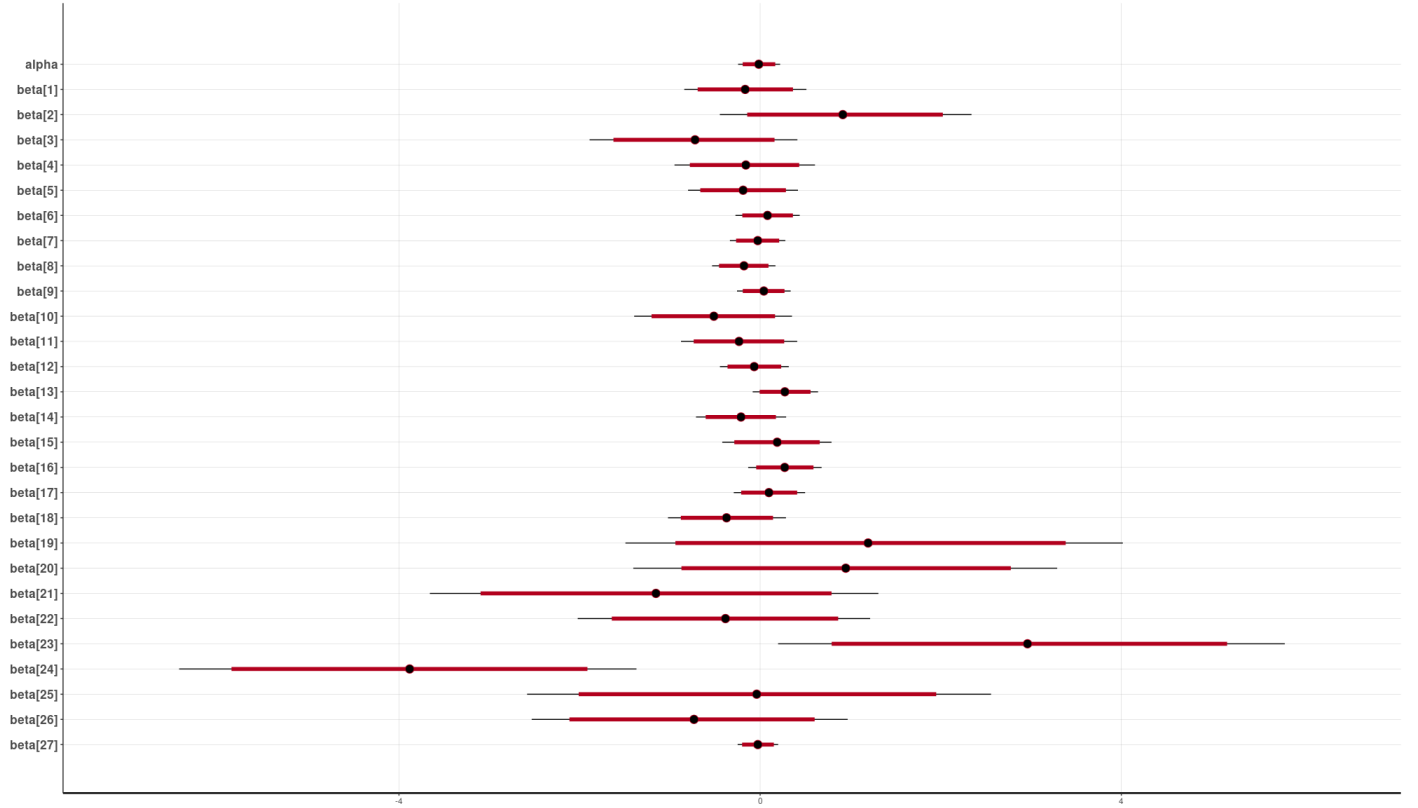


Figure 3: The 80% and 90% credible intervals for the parameters of our best model, based on the approach of Model 2. Barring the intercept term α , the order in which the parameters are listed corresponds to the order in which the variables are listed in table 1

6 Discussion

We first note that the cross validation accuracy results listed in table 3 indicate that our model is not really sensitive to prior specification. Specifically, we observe that the cross validation accuracy for Model 2 and Model 3 only differ by 0.4%. Model 2 only slightly outperforms Model 3. Furthermore, the results in Section 5 indicate that the most predictive covariates for the winner of a March Madness game are the difference of their average Pomeroy and B Wilson Empirical rankings. In particular, holding all other predictors constant, the expected change in the log odds ratio of winning a game when the POMMeanDiff between the first listed team and the second listed team increases by one standard deviation is around -4. Similarly, holding all other predictors constant, the expected change in the log odds ratio of winning a game when the BWEMeanDiff between the first listed team and the second listed team increases by one standard deviation is around positive 3. We note that the coefficient for POMMeanDiff makes sense; as the difference in the average Pomeroy ranking between the first team and the second team becomes more positive i.e the ordinal rank of the first team “decreases” (i.e increases in magnitude) or the ordinal rank of the second team “increases” (i.e decreases in magnitude) - and hence the difference between their average rankings becomes more positive - the log of the odds of winning decreases. On the other hand, the coefficient for BWEMeanDiff does not make a lot sense,

because its interpretation is the opposite of that of the average Pomeroy difference coefficient. Given the high quality of the BWE ranking system, we suspect that it our model that is flawed, rather than the ranking system being not credible. This is especially true given the mediocre accuracy results obtained by our model.

7 Conclusion

There were a plethora of challenges associated with our work which may have contributed to the mediocre results we obtained, especially with our test set. There are over 350 teams in Division I college basketball and only 68 make the tournament each year. On top of that, the rosters for each team vary greatly for each year, making year to year trends hard to predict. Additionally, college student athletes are amateurs and not professionals, making it hard to predict consistency. Unfortunately, the 2022 tournament was even more unpredictable than usual since there were many unexpected upsets including a 15 seed making the quarter finals for the first time, an 8 seed making the championship for the third time, and a few other double-digit seeds making it past the second round. It may be that logistic regression is not powerful enough to model the “madness” of March Madness. As an extension to the current work, we suggest testing more powerful models like Random Forest or XGBoost. Overall, we determined that it is very difficult to predict the outcomes of games in the NCAA tournament. Perfect accuracy may be nearly impossible to obtain. Indeed, Warren Buffett offers \$1 billion to the first person who can accurately predict each game in the tournament, hinting at the near impossibility of the task given current statistical methods.

References

- Brown, B. (2019). Predictive analytics for college basketball: Using logistic regression for determining the outcome of a game regression for determining the outcome of a game. <https://scholars.unh.edu/cgi/viewcontent.cgi?article=1421&context=honors>
- Nelson, B. (2012). Modeling the ncaa tournament through bayesian logistic regression. <https://dsc.duq.edu/cgi/viewcontent.cgi?article=1986&context=etd>

8 Appendix

8.1 Header

```
1 library(plyr)
2 library(tidyverse)
3 library(rstan)
4 options(max.print = 100000)
5 options(mc.cores = parallel::detectCores())
6 set.seed(42)
```

8.2 Data Prep

```
1 ##### Data Prep #####
2 marchMadnessAllGames <- read_csv(
3   "MDataFiles_Stage2/MNCAATourneyCompactResults.csv")
4 regSeason <- read_csv("MDataFiles_Stage2/MRegularSeasonDetailedResults.csv")
5 seeds <- read_csv("MDataFiles_Stage2/MNCAATourneySeeds.csv")
6 ranks <- read_csv("MDataFiles_Stage2/MMasseyOrdinals_thruDay128.csv")
7 teamIDToTeamName <- read_csv("MDataFiles_Stage2/MTeams.csv")
8
9 # Specify list of desired seasons for teams
10 # returns MM teams participating in MM seasons specified by "seasons"
11 getMMTeams <- function(seasons, marchMadnessGames) {
12   marchMadnessGames <- marchMadnessGames %>%
13     filter(Season %in% seasons)
14   wteams <- marchMadnessGames %>%
15     select(WTeamID)
16   lteams <- marchMadnessGames %>%
17     select(LTeamID)
18   marchMadnessTeams <- c(unlist(wteams$WTeamID), unlist(lteams$LTeamID))
19   marchMadnessTeams <- unique(marchMadnessTeams)
20   return(marchMadnessTeams)
21 }
22
23 # Must pass in all seasons of interest
24 getAvailableRankingSystems <- function(seasons, ranks, marchMadnessGames) {
25   rankingSystems <- unique(ranks$SystemName)
26   marchMadnessTeams <- getMMTeams(seasons, marchMadnessGames)
27   ranks <- ranks %>%
28     filter(Season %in% seasons, TeamID %in% marchMadnessTeams)
29
30   availableRankingsAggregated <- list()
31   for (team in marchMadnessTeams) {
32     for (season in unique(ranks$Season)) {
33       availableRankings <- ranks %>%
34         filter(TeamID == team, Season == season) %>%
35         select(SystemName) %>%
36         unique()
37       availableRankingsAggregated <- availableRankingsAggregated %>%
38         append(availableRankings)
39     }
40   }
41   return(Reduce(intersect, availableRankingsAggregated))
42 }
43
44 # UNCOMMENT and run if want to check for different seasons
45 # availableRankingSystems <- getAvailableRankingSystems(c(2016:2022),
46 #   ranks,
47 #   marchMadnessAllGames)
48 # availableRankingSystems <- getAvailableRankingSystems(c(2003:2022),
49 #   ranks,
50 #   marchMadnessAllGames)
51
52 getSeedsPerSeason <- function(season, seeds) {
53   seeds <- seeds %>%
54     filter(Season == season)
55   seeds <- seeds %>%
56     mutate(Seed = str_sub(Seed, 2, 3))
57   seeds$Seed <- as.integer(seeds$Seed)
58   return(seeds)
59 }
60
61 # returns average season rank in each desired ranking system and last ranking
62 # right before beginning of March Madness for each team (rankingDayNum 133
63 # is not always available before the start of the tournament, so we use Day 128
64 # instead)
65 getRanksPerSeason <- function(season, desiredRankingSystems, ranks,
66   teams) {
67   preMMranks <- ranks %>%
68     filter(Season == season, SystemName %in% desiredRankingSystems,
69           TeamID %in% teams, RankingDayNum == 128) %>%
```

```

70     select(-one_of("Season", "RankingDayNum")) %>%
71     pivot_wider(names_from = SystemName, values_from = OrdinalRank)
72
73     averageRegSeasonranks <- ranks %>%
74       filter(Season == season, SystemName %in%
75         desiredRankingSystems, TeamID %in% teams) %>%
76       group_by(TeamID, SystemName) %>%
77       summarize(mean(OddinalRank)) %>%
78       pivot_wider(names_from = SystemName, values_from =
79         "mean(OddinalRank)")
80   if (nrow(averageRegSeasonranks) * ncol(averageRegSeasonranks) != 0) {
81     colnames(averageRegSeasonranks) <- c("TeamID", paste(desiredRankingSystems,
82       "Mean", sep=''))
83   }
84   ranks <- preMMranks %>%
85     inner_join(averageRegSeasonranks, by="TeamID")
86   return(ranks)
87 }
88
89 getMMTeamsPreTourneyStatsPerSeason <- function(season, regSeason,
90   marchMadnessAllGames,
91   seeds, ranks,
92   desiredRankingSystems) {
93   regSeason <- regSeason %>%
94     filter(Season == season)
95   seeds <- getSeedsPerSeason(season, seeds)
96   MMTeams <- getMMTeams(season, marchMadnessAllGames)
97   ranks <- getRanksPerSeason(season, desiredRankingSystems, ranks, MMTeams)
98   regSeason <- regSeason %>%
99     filter(WTeamID %in% MMTeams | LTeamID %in% MMTeams)
100   dup1 <- regSeason %>%
101     mutate(TeamID=WTeamID)
102   dup2 <- regSeason %>%
103     mutate(TeamID=LTeamID)
104   regSeason <- bind_rows(dup1, dup2)
105   MMTeamsRegSeasonStats <- tibble()
106   for (team in MMTeams) {
107     teamTibble <- regSeason %>%
108       filter(TeamID == team)
109     winTibble <- teamTibble %>%
110       filter(WTeamID == team)
111     winTibble <- winTibble %>%
112       select(-one_of("DayNum", "WPF", "LPF", "NumOT", "WLoc",
113         "TeamID"))
114     colnames(winTibble) <- c("Season", "currTeamID", "currScore", "otherTeamID",
115       "otherScore", "currFGM", "currFGA", "currFGM3",
116       "currFGA3", "currFTM", "currFTA", "currOR",
117       "currDR", "currAst", "currTO", "currStl",
118       "currBlk", "otherFGM", "otherFGA", "otherFGM3",
119       "otherFGA3", "otherFTM", "otherFTA", "otherOR",
120       "otherDR", "otherAst", "otherTO", "otherStl",
121       "otherBlk")
122     loseTibble <- teamTibble %>%
123       filter(LTeamID == team)
124     loseTibble <- loseTibble %>%
125       select(-one_of("DayNum", "WPF", "LPF", "NumOT", "WLoc",
126         "TeamID"))
127     colnames(loseTibble) <- c("Season", "otherTeamID", "otherScore",
128       "currTeamID", "currScore", "otherFGM", "otherFGA",
129       "otherFGM3", "otherFGA3", "otherFTM", "otherFTA",
130       "otherOR", "otherDR", "otherAst", "otherTO",
131       "otherStl", "otherBlk", "currFGM", "currFGA",
132       "currFGM3", "currFGA3", "currFTM", "currFTA",
133       "currOR", "currDR", "currAst", "currTO", "currStl",
134       "currBlk")
135     currTeamTibble <- bind_rows(winTibble, loseTibble)
136     currTeamTibble <- currTeamTibble %>%
137       summarize(TeamID=unique(currTeamID),
138         winPct = nrow(winTibble) / (nrow(winTibble) +
139           nrow(loseTibble)),
140         lossPct = nrow(loseTibble) / (nrow(winTibble) +
141           nrow(loseTibble)),
142         currScoreAvg = mean(currScore),
143         otherScoreAvg = mean(otherScore),
144         currFGPct = mean(currFGM)/mean(currFGA),
145         otherFGPct = mean(otherFGM)/mean(otherFGA),
146         currFG3Pct = mean(currFGM3)/mean(currFGA3),
147         otherFG3Pct = mean(otherFGM3)/mean(otherFGA3),
148         currFTPct = mean(currFTM)/mean(currFTA),
149         otherFTPct = mean(otherFTM)/mean(otherFTA),
150         currRebAvg = (mean(currOR) + mean(currDR)),
151         otherRebAvg = (mean(otherOR) + mean(otherDR)),
152         currAstAvg = mean(currAst),
153         otherAstAvg = mean(otherAst),
154         currTOAvg = mean(currTO),
155         otherTOAvg = mean(otherTO),
156         currDefAvg = (mean(currStl) + mean(currBlk)),
157         otherDefAvg = (mean(otherStl) + mean(otherBlk))

```

```

158     )
159     MMTeamsRegSeasonStats <- MMTeamsRegSeasonStats %>%
160       bind_rows(currTeamTibble)
161   }
162
163   MMTeamsRegSeasonStats <- MMTeamsRegSeasonStats %>%
164     inner_join(seeds, by="TeamID")
165
166   if (nrow(ranks) * ncol(ranks) != 0) {
167     MMTeamsRegSeasonStats <- MMTeamsRegSeasonStats %>%
168       inner_join(ranks, by="TeamID")
169   }
170   return(MMTeamsRegSeasonStats)
171 }
172
173 computeStatDiff <- function(df1, df2) {
174   diff <- df1 - df2
175   colnames(diff) <- paste(colnames(diff), "Diff", sep="")
176   return(diff)
177 }
178
179 computeSelectStatsAddition <- function(df1, df2, selectColumns) {
180   df1 <- df1 %>%
181     select(selectColumns)
182
183   df2 <- df2 %>%
184     select(selectColumns)
185   addition <- df1 + df2
186   colnames(addition) <- paste(colnames(addition), "Add", sep="")
187   return(addition)
188 }
189
190 # season must match regSeasonStats season
191 getMarchMadnessMatchupsPerSeason <- function(season, MMTeamsRegSeasonStats,
192                                               marchMadnessAllGames) {
193   MMTeamsRegSeasonStats <- MMTeamsRegSeasonStats %>% select(-Season)
194   marchMadnessGames <- marchMadnessAllGames %>%
195     filter(Season == season) %>%
196     sample_frac(1L)
197   num_row <- nrow(marchMadnessGames)
198   marchMadnessWins <- marchMadnessGames %>%
199     slice(1:floor(num_row/2)) %>%
200     select(WTeamID, LTeamID) %>%
201     mutate(firstTeamWinStatus = 1)
202   colnames(marchMadnessWins)[1] <- "TeamID"
203   marchMadnessLosses <- marchMadnessGames %>%
204     slice(floor(num_row/2)+1:num_row) %>%
205     select(LTeamID, WTeamID) %>%
206     mutate(firstTeamWinStatus = 0)
207   colnames(marchMadnessLosses)[1] <- "TeamID"
208
209   getTeamStatsByKey <- function(df, regSeasonStats, Key) {
210     marchMadnessFirstTeamStats <- df %>%
211       left_join(regSeasonStats,
212               by=Key) %>%
213       select(-c(1,2,3))
214     return(marchMadnessFirstTeamStats)
215   }
216
217   marchMadnessWinsWinTeamStats <- getTeamStatsByKey(marchMadnessWins,
218                                                     MMTeamsRegSeasonStats,
219                                                     "TeamID")
220   marchMadnessWinsLoseTeamStats <- getTeamStatsByKey(marchMadnessWins,
221                                                       MMTeamsRegSeasonStats,
222                                                       c("LTeamID" = "TeamID"))
223   marchMadnessWinsStatDiff <- computeStatDiff(marchMadnessWinsWinTeamStats,
224                                                marchMadnessWinsLoseTeamStats) %>%
225     select(-lossPctDiff)
226   marchMadnessWinsStatAdd <- computeSelectStatsAddition(
227     marchMadnessWinsWinTeamStats,
228     marchMadnessWinsLoseTeamStats,
229     c("Seed"))
230
231   marchMadnessWins <- marchMadnessWins %>%
232     bind_cols(marchMadnessWinsStatDiff) %>%
233     bind_cols(marchMadnessWinsStatAdd)
234
235   marchMadnessLossesLoseTeamStats <- getTeamStatsByKey(marchMadnessLosses,
236                                                         MMTeamsRegSeasonStats,
237                                                         "TeamID")
238   marchMadnessLossesWinTeamStats <- getTeamStatsByKey(marchMadnessLosses,
239                                                         MMTeamsRegSeasonStats,
240                                                         c("WTeamID" = "TeamID"))
241   marchMadnessLossesStatDiff <- computeStatDiff(marchMadnessLossesLoseTeamStats,
242                                                  marchMadnessLossesWinTeamStats) %>%
243     select(-lossPctDiff)
244   marchMadnessLossesStatAdd <- computeSelectStatsAddition(
245     marchMadnessLossesLoseTeamStats,

```

```

246                                     marchMadnessLossesWinTeamStats,
247                                     c("Seed"))
248 marchMadnessLosses <- marchMadnessLosses %>%
249   bind_cols(marchMadnessLossesStatDiff) %>%
250   bind_cols(marchMadnessLossesStatAdd)
251 colnames(marchMadnessLosses)[2] <- "otherTeamID"
252 colnames(marchMadnessWins)[2] <- "otherTeamID"
253 clean_data <- marchMadnessLosses %>%
254   bind_rows(marchMadnessWins) %>%
255   sample_frac(1L)
256
257 Season <- as_tibble(rep(season, nrow(clean_data)))
258 colnames(Season) <- "Season"
259 clean_data <- bind_cols(Season, clean_data)
260 return(clean_data)
261 }
262
263 createDataset <- function(seasons, regSeason, marchMadnessAllGames, seeds,
264   ranks, desiredRankingSystems, teamIDToTeamName) {
265   MMTTeamsPreTourneyStats <- list()
266   MMMatchups <- list()
267   for (i in 1:length(seasons)) {
268     season <- seasons[i]
269     MMTTeamsPreTourneyStats <- append(MMTTeamsPreTourneyStats,
270   list(getMMTeamsPreTourneyStatsPerSeason(
271     season,
272     regSeason,
273     marchMadnessAllGames,
274     seeds,
275     ranks,
276     desiredRankingSystems)))
277     MMTTeamsPreTourneyStats[[i]] <- MMTTeamsPreTourneyStats[[i]] %>%
278   select(Season, everything())
279     MMMatchups <- append(MMMatchups, list(getMarchMadnessMatchupsPerSeason(
280   season,
281   MMTTeamsPreTourneyStats[[i]],
282   marchMadnessAllGames)))
283   }
284   MMDataset <- do.call(rbind, MMMatchups)
285
286   MMDataset$TeamID <- (MMDataset %>% left_join(
287     teamIDToTeamName, on="TeamID"))$TeamName
288
289   colnames(MMDataset)[2] <- "firstTeamName"
290   colnames(MMDataset)[3] <- "otherTeamName"
291
292   MMDataset$otherTeamName <- (MMDataset %>% left_join(
293     teamIDToTeamName, by=c("otherTeamName" = "TeamID")))$TeamName
294
295   return(MMDataset)
296 }
297
298 MM2016To2021Dataset <- createDataset(c(2016, 2017, 2018, 2019, 2021), regSeason,
299   marchMadnessAllGames, seeds, ranks,
300   c("POM", "BWE", "MAS", "SFX"),
301   teamIDToTeamName)
302 # Descriptive statistics
303 MM2016To2021DatasetMeans <- MM2016To2021Dataset %>%
304   select(-one_of("Season", "firstTeamName",
305     "otherTeamName",
306     "firstTeamWinStatus")) %>%
307   summarise(across(everything(), mean))
308
309 MM2016To2021DatasetSds <- MM2016To2021Dataset %>%
310   select(-one_of("Season", "firstTeamName",
311     "otherTeamName",
312     "firstTeamWinStatus")) %>%
313   summarise(across(everything(), sd))
314
315 seasons <- c(2003:2021)
316 seasons <- seasons[seasons != 2020]
317 MM2003To2021Dataset <- createDataset(seasons, regSeason,
318   marchMadnessAllGames, seeds, ranks,
319   c(), teamIDToTeamName)

```

8.3 Stan Modeling

```

1 ##### Modeling #####
2 createModel <- function() {
3   if (!file.exists("logisticRegressionModel.stan")) {
4     write("
5     data {
6       int<lower=0> N; // number of data items
7       int<lower=0> K; // number of predictors
8       matrix[N, K] X; // predictor matrix

```

```

9       int y[N];           // outcome vector
10      int relativePrior;
11      int shrinkageIdx[4];
12      int noShrinkageIdx[23];
13    }
14    parameters {
15      real alpha;           // intercept
16      vector[K] beta;       // coefficients for predictors
17    }
18    model {
19      // Priors:
20      alpha ~ normal(0, 10);
21      if (relativePrior) {
22        for (k in shrinkageIdx) {
23          beta[k] ~ normal(0, 1);
24        }
25        for (k in noShrinkageIdx) {
26          beta[k] ~ normal(0, 10);
27        }
28      }
29      else {
30        beta ~ normal(0, 10);
31      }
32      // Likelihood
33      y ~ bernoulli_logit(alpha + X * beta);
34    }
35    generated quantities {
36      vector[N] y_preds;
37      real correct = 0;
38      real accuracy;
39      for (n in 1:N) {
40        y_preds[n] = bernoulli_logit_rng(alpha + X[n] * beta);
41        correct += logical_eq(y_preds[n], y[n]);
42      }
43      accuracy = correct / N;
44    }
45    ", "logisticRegressionModel.stan")
46  }
47 }
48
49 createModel()

```

8.4 Diagnostics

8.4.1 Cross Validation

```

1  ##### Diagnostics #####
2  scaleTrainData <- function(data, exclude_columns) {
3    covariates <- data %>% select(-one_of(c(unlist(exclude_columns))))
4    means <- c()
5    stdevs <- c()
6    for (column in colnames(covariates)) {
7      mean <- mean(covariates[[column]])
8      means <- append(means, mean)
9      stdev <- sd(covariates[[column]])
10     stdevs <- append(stdevs, stdev)
11     covariates[column] <- scale(covariates[column])
12   }
13   return(list(covariates, means, stdevs))
14 }
15
16 scaleTestData <- function(data, means, stdevs, exclude_columns) {
17   covariates <- data %>% select(-one_of(c(unlist(exclude_columns))))
18   for (i in 1:length(colnames(covariates))) {
19     column <- colnames(covariates)[i]
20     covariates[column] <- (covariates[column] - means[i]) / stdevs[i]
21   }
22   return(covariates)
23 }
24 createStanDataLists <- function(dataset, val_fold, relativePriors,
25                                outcome_column_name) {
26   outcome <- dataset[outcome_column_name]
27   outcome_train <- (dataset[dataset$fold != val_fold,][outcome_column_name])
28   outcome_val <- (dataset[dataset$fold == val_fold,][outcome_column_name])
29   dataset_train <- dataset[dataset$fold != val_fold,]
30   dataset_val <- dataset[dataset$fold == val_fold,]
31   exclude_columns <- c("firstTeamName", "otherTeamName", outcome_column_name,
32                        "fold")
33   scaleList <- scaleTrainData(dataset_train, exclude_columns)
34   covariates_train <- scaleList[[1]]
35   means <- scaleList[[2]]
36   stdevs <- scaleList[[3]]
37   covariates_val <- scaleTestData(dataset_val, means, stdevs, exclude_columns)
38   noShrinkageIdx <- c(1:27)
39   shrinkageIdx <- c(20, 22, 25, 26)

```

```

40 noShrinkageIdx <- noShrinkageIdx[!(noShrinkageIdx %in% shrinkageIdx)]
41 train_list <- list(X=covariates_train, y=array(unlist(outcome_train)),
42                  K=ncol(covariates_train),
43                  N=nrow(covariates_train),
44                  relativePrior=relativePriors,
45                  shrinkageIdx=shrinkageIdx,
46                  noShrinkageIdx=noShrinkageIdx
47                  )
48
49 val_list <- list(X=covariates_val, y=array(unlist(outcome_val)),
50                K=ncol(covariates_val),
51                N=nrow(covariates_val),
52                relativePrior=relativePriors,
53                shrinkageIdx=shrinkageIdx,
54                noShrinkageIdx=noShrinkageIdx
55                )
56 return(list(train_list, val_list))
57 }
58
59 CV <- function(dataset, relativePriors) {
60   if (length(unique(dataset$Season)) == 5) {
61     dataset$fold <- dataset$Season - 2016
62     dataset$fold[dataset$fold == 5] <- 4
63   }
64   else {
65     dataset$fold <- mapvalues(dataset$Season,
66                             from=c("2003", "2004", "2005", "2006", "2007", "2008",
67                                   "2009", "2010", "2011", "2012", "2013", "2014",
68                                   "2015", "2016", "2017", "2018", "2019", "2021"),
69                             to=c(rep(0, 3), rep(1, 3),
70                                rep(2, 3), rep(3, 3),
71                                rep(4, 3), rep(5, 3)))
72   }
73   dataset <- dataset %>% select(-Season)
74   to_return <- list()
75   fold_fits <- list()
76   valPerformances <- list()
77   for(fold in 0:(length(unique(dataset$fold))-1)) {
78     stanDataLists <- createStanDataLists(dataset, fold, relativePriors,
79                                         "firstTeamWinStatus")
80     trainList <- stanDataLists[[1]]
81     valList <- stanDataLists[[2]]
82     fit <- stan(file = "logisticRegressionModel.stan", data = trainList,
83                iter=10000, seed=42)
84     fold_fits <- append(fold_fits, fit)
85     valPerformances <- append(valPerformances, gqs(fit@stanmodel,
86                                                    draws = as.matrix(fit),
87                                                    data = valList))
88   }
89   to_return <- list(fold_fits, valPerformances)
90   return(to_return)
91 }
92
93 fits20162021NoRelativeRankingPrior <- CV(MM2016To2021Dataset, 0)
94 avgAccuracy20162021NoRelativeRankingPrior <- mean(c(66, 67, 66, 66, 67))
95 fits20162021WithRelativeRankingPrior <- CV(MM2016To2021Dataset, 1)
96 avgAccuracy20162021WithRelativeRankingPrior <- mean(c(66, 66, 66, 65, 67))
97 fits20032021NoRanking <- CV(MM2003To2021Dataset, 0)
98 avgAccuracy20032021NoRanking <- mean(c(63, 63, 64, 64, 64))

```

8.4.2 Train best model

```

1 # Train best model
2 MM2016To2021DatasetDup <- MM2016To2021Dataset
3 MM2016To2021DatasetDup$fold <- 0
4 MM2016To2021DatasetDup <- MM2016To2021DatasetDup %>% select(-Season)
5 outcome <- MM2016To2021DatasetDup["firstTeamWinStatus"]
6 exclude_columns <- c("firstTeamName", "otherTeamName", "firstTeamWinStatus",
7                      "fold")
8 scaleList <- scaleTrainData(MM2016To2021DatasetDup, exclude_columns)
9 covariates <- scaleList[[1]]
10 means <- scaleList[[2]]
11 stdevs <- scaleList[[3]]
12 noShrinkageIdx <- c(1:27)
13 shrinkageIdx <- c(20, 22, 25, 26)
14 noShrinkageIdx <- noShrinkageIdx[!(noShrinkageIdx %in% shrinkageIdx)]
15 train_list <- list(X=covariates, y=array(unlist(outcome)),
16                  K=ncol(covariates),
17                  N=nrow(covariates),
18                  relativePrior=0,
19                  shrinkageIdx=shrinkageIdx,
20                  noShrinkageIdx=noShrinkageIdx
21                  )
22 fit <- stan(file = "logisticRegressionModel.stan", data = train_list,
23            iter=10000, seed=42)

```



```

24 pars <- c('alpha', 'beta[1]', 'beta[2]', 'beta[3]', 'beta[4]', 'beta[5]',
25           'beta[6]', 'beta[7]', 'beta[8]', 'beta[9]', 'beta[10]', 'beta[11]',
26           'beta[12]', 'beta[13]', 'beta[14]', 'beta[15]', 'beta[16]', 'beta[17]',
27           'beta[18]', 'beta[19]', 'beta[20]', 'beta[21]', 'beta[22]',
28           'beta[23]', 'beta[24]', 'beta[25]', 'beta[26]', 'beta[27]')
29
30 traceplot(fit, pars=pars)
31 covariates <- colnames(MM2016To2021DatasetDup)[4:30]

```

8.4.3 Test set Performance

```

1  # Model Test Set Performance
2  best_model <- fit
3  set.seed(42)
4  testSet <- createDataset(2022, regSeason,
5                           marchMadnessAllGames, seeds, ranks,
6                           c("POM", "BWE", "MAS", "SFX"), teamIDToTeamName)
7  testSet <- testSet %>% select(-Season)
8  outcome <- testSet["firstTeamWinStatus"]
9  testSet <- scaleTestData(testSet, means, stdevs, exclude_columns)
10 testList <- list(X=testSet, y=array(unlist(outcome)),
11                 K=ncol(testSet),
12                 N=nrow(testSet),
13                 relativePrior=0,
14                 shrinkageIdx=shrinkageIdx,
15                 noShrinkageIdx=noShrinkageIdx
16 )
17 gqs(fit@stanmodel, draws=as.matrix(best_model), data=testList, seed=42)

```

8.5 Model Interpretation

```

1  #### Model interpretation ####
2  plot(fit, pars=pars, outer_level=0.90)

```