# NCAA Men's Basketball Rankings: Reconstructing

# of the Team Value Index

**By Md. Rifayat Uddin**

**Preface (Cover Letter)**

After getting very good feedback from the Instructor and my colleagues, I came to the conclusion that my paper needed minor improvements in the introduction, and a few major improvements in the methodology and results section. The methodology section, in particular, needed the most revision. This is primarily because I was ignoring taking into account the multiplicity when conducting the hypothesis tests. I was also given the feedback that some of the images and tables aren't as interpretable as it would've been for a student studying Data Science. Hence, I decided to make sure I simplified these tables for ease of the reader.

In order to improve my paper, I got rid of unnecessary information, especially portions that used the "glm" and "brglm" functions instead of explicitly stating what they are. The only suggestion I didn't end up using was adding results in my introduction since there would be a lot of things to include. Hence, I decided to keep that for the results section. Overall, this was a very enjoyable paper for me to write on.

**Introduction**

NCAA Men's Basketball Tournament is one of the most followed and watched college athletics tournaments in the US. In fact, last year in 2018, over 97 million viewers in 180 countries watched the first four to final four leg of the tournament.(1) The NCAA website reports that the national championship game attendance has gone up from a little under 10,000 people attending back I 1939 to over 70,000 people attending the final fame last year. The interest of people in the NCAA games is high, which is why rankings between these college teams play a big role in designing brackets for March Madness and helps create excitement for the games overall.

The rankings of the teams are usually done by taking five components into account; one of which isn't publicly available called the Team Value Index (TVI). This has led to speculation within statisticians to approximate this TVI through statistical models. This paper addresses this TVI through the creation of a proxy for the undisclosed TVI portion of the ranking's calculation. This is done using the Bradley-Terry models(2) to fit models taking into account the factors such as game locations to get the best possible model. The data used for this paper was sourced from the game results data available through CBS Sports for the time-frame until mid-January. The paper addresses two key questions; whether the home court advantage (or game location for that matter) plays significant role in rankings and whether Michigan Men's Basketball team is inferior to Illinois, OSU and MSU's Men's Basketball teams.

**Methodology**

The paper first uses the Bradley-Terry model on the data by regressing over whether or not the home team won based on who the home team was and against whom they were playing using using two different fitting methods; standard fitting and a fitting method that includes bias reduction in its fitting to generate the rankings. We ended up choosing the fitting method that includes bias reduction since it helps address the bias. This was also complemented by the fact that the rankings didn't change while using the two different functions. Since the data is relatively small, with only 29 observations, it made calculations extremely fast.

Furthermore, to test for whether the team rankings were accurately being represented by the TVI, a smaller test was done on the data whereby through a function, the team ability index was tested to predict which team won in a game for all the teams in the models except Rutgers. The test counted the number of games each team would win and based on that, a ranking was created. This made model fitting simple and easier to run. Upon conducting this test, it was found that the rankings didn't differ from the one created by the Bradley-Terry model.

To address the first research question of whether home-court advantage and game location plays a significant role in rankings, a model was fit including the intercept that is considered the game-location coefficient in the model. Then, the full model and the home-court advantage models were compared using ANOVA to test for any significant difference. The rankings of teams for both the models was compared to check for any reportable differences. Game location can be an important factor in determining the TVI for the teams, as familiarity

with a location can help the team play better in certain locations versus others. This is why testing for whether game location is an important factor to look at was extremely important.

For the second research question, the paper looks at whether Michigan Men's Basketball team is inferior in comparison to three other teams of the author's choice – in this paper's case, these teams are Illinois, Ohio State and Michigan State Basketball teams. In order to answer this research question, a hypothesis test was done, called the "non-inferiority" test. To answer this question, a "non-inferiority" test was conducted that tests three hypothesis tests simultaneously with the hypothesis of "true ability coefficient of Michigan – true ability coefficient of the team to be tested on is less than or equal to 0". This question is also important for me personally, as it helps understand where Michigan stands statistically in comparison to these very popular teams in the Big Ten.

Finally, the paper also looks at comparing a model generated by fitting a Bayesian Generalized Linear model with the chosen model between the general model and the model that accounts for the home-court advantage. The comparison was done by looking at cross-validation errors and finding out which model is the best fit for this data.

**Results**

Figure 1 summarizes the team ability indices for each team found from fitting the Bradley-Terry models using the two methods – Standard Bradley Terry and Bradley Terry with Bias Reduction and contains the team ability indices produced by the two methods as well.

***Figure 1: Rankings and Model Information Comparing Standard Bradley Terry and Bradley Terry with Bias Reduction***

| | Bradley Terry | Bradley Terry with Bias Reduction |
|---|---|---|
| Degrees of Freedom | 16 | 16 |
| Residual Deviance | 40.2 | 22.7 |
| University of Michigan Rank | 1 | 1 |
| Purdue University Rank | 2 | 2 |
| Michigan State University Rank | 3 | 3 |
| University of Maryland Rank | 4 | 4 |
| Indiana University Rank | 5 | 5 |
| University of Wisconsin Rank | 6 | 6 |
| University of Minnesota Rank | 7 | 7 |
| Ohio State University Rank | 8 | 8 |
| University of Iowa Rank | 9 | 9 |
| Penn State University Rank | 10 | 10 |
| Northwestern University Rank | 11 | 11 |
| University of Nebraska Rank | 12 | 12 |
| University of Illinois Rank | 13 | 13 |

*The Figure above shows the team ability indices produced by the two methods showing the rankings that are the same. While neither estimate is seen to be statistically significant, the p-values for Bias reduction method are significantly lower.*

As mentioned earlier, the fitting method using the bias is used for the paper as the p-values of the estimates are also significantly lower than the standard fitting method.

Furthermore, Table 2 in the next page shows the rankings for teams generated from the Bradley-Terry model in the first column and the rankings predicted by the test that was run on the data. As is evident, the Michigan Basketball Men's Team is ranked higher than the rest of the teams with 13 wins and Nebraska and Illinois in the bottom of the rankings according to the test. However, the model predicts Illinois to have a lower ranking compared to Nebraska as it takes intricacies of the data and analyses it better than just number of wins.

***Table 2: Rankings(Descending Order)***

*The table shows that the rankings for both methods are the same. The Number of possible wins however is a least of 1 for both Nebraska and Illinois which places them in rank12 for the test method, but the model method accurately places Nebraska over Illinois*

|  | Rankings predicted by Model (for both methods) | Ranking Predicted by Number of possible wins (test) | Number of wins (predicted by test) |
|---|---|---|---|
| 1 | University of Michigan | University of Michigan | 12 |
| 2 | Purdue University | Purdue University | 11 |
| 3 | Michigan State University | Michigan State University | 10 |
| 4 | University of Maryland | University of Maryland | 9 |
| 5 | Indiana University | Indiana University | 8 |
| 6 | University of Wisconsin | University of Wisconsin | 7 |
| 7 | University of Minnesota |  | 6 |
| 8 | Ohio State University | Ohio State University | 5 |
| 9 | University of Iowa | University of Iowa | 4 |
| 10 | Penn State University | Penn State University | 3 |
| 11 | Northwestern University | Northwestern University | 2 |
| 12 | University of Nebraska | University of Nebraska + University of Illinois | 1 |
| 13 | University of Illinois |  |  |

In order to fit the model, the intercept term which is considered to be the location coefficient was taken out. However, to test for whether game location is an significant factor to consider, the analysis of variance test (ANOVA) was conducted by testing on a model with the game location coefficient and without it. Since the p-value was not found, it's assumed that the models are not significantly different from one another. Upon further analysis, it was also found that the rankings did change. Since the model without the game location as a factor has lower residual deviance, and higher degrees of freedom, it can be concluded that the game location is not a statistically significant factor in determining the TVI for the teams. Table 3 below summarizes the rankings for the teams generated by the model with game location and model without.

Table 3: **Rankings without Game Location vs Rankings with Game Location**

*Table 3 shows that if game location is taken into account, teams in the red colored cell see a drop in ranking and teams in green colored cell see rise. These teams corresponding locations in the table can be found in the right where the font is shaded in their corresponding colors.*

| | Rankings predicted by Model | Rankings predicted by Model + Game Location |
|---|---|---|
| 1 | University of Michigan | University of Michigan |
| 2 | Purdue University | Purdue University |
| 3 | Michigan State University | Michigan State University |
| 4 | University of Maryland | University of Maryland |
| 5 | Indiana University | **University of Minnesota** |
| 6 | University of Wisconsin | University of Wisconsin |
| 7 | University of Minnesota | **Indiana University** |
| 8 | Ohio State University | Ohio State University |
| 9 | University of Iowa | University of Iowa |
| 10 | Penn State University | **Northwestern University** |
| 11 | Northwestern University | **Penn State University** |
| 12 | University of Nebraska | University of Nebraska |
| 13 | University of Illinois | University of Illinois |

Based on the statistical comparison on both the models, and since only 4 teams see a switch in rankings by a maximum of 2 spots, game location is concluded to not be a significant factor in determining the TVI of the teams.

In order to answer the second research question for the paper, which is whether Michigan is inferior to the teams of Illinois, Michigan State and Ohio State, a hypothesis test was conducted. The hypothesis test used was the non-inferiority test where the null-hypothesis uses the difference between the estimated TVI's of Michigan and the compared teams to be less than or equal to zero. The table below summarizes the p-values found from conducting the hypothesis tests and the decision for each test. Upon conducting the test, it is found that for the test between Michigan and Illinois, the p-value found is not less than the significance level of 0.05 so we fail to reject the null-hypothesis and cannot conclude whether Michigan is inferior to Illinois. The same results are observed for the tests between Michigan and Ohio State and

Michigan and Penn State. Thus, we cannot conclude anything about the inferiority of Michigan in comparison to these teams. It's important to note that we use the Bonferroni-Holm method to account for the multiplicity of the model.

**Table 4: Hypothesis Test Outcomes for Inferiority Tests**
*The table below summarizes the hypothesis test that Michigan is not inferior to Illinois, OSU and PSU and the p-values show we fail to reject the null-hypothesis*

| Game | Adjusted P-Value | Hypothesis Decision |
|---|---|---|
| Michigan Vs Illinois | 0.1632 | Fail to Reject |
| Michigan Vs OSU | 0.3142 | Fail to Reject |
| Michigan Vs Penn State | 0.2963 | Fail to Reject |

Furthermore, the paper also addresses which model is the best one to choose from among the Bradley-Terry model, Bradley-Terry model with Bias Reduction and a Bayesian model. The table 5 below summarizes the results found using cross-validation errors by dividing the data 10-times and running it over 10-times to improve accuracy.  As the table summarizes, we will end up choosing the Bradley-Terry model instead since it has the lowest cross-validation error.

**Table 5: Cross-Validation Error table for Bradley-Terry, Bradley Terry with Bias Reduction and Bayesian Model**
*The Table below shows that the Cross-validation error for the Bradley-Terry model is the lowest, hence statistically, we should select the Bradley-Terry model*

| Model | Cross-Validation Error |
|---|---|
| Bayesian Model | 0.23 |

| Bradley-Terry Model with Bias Reduction | 0.25 |
| Bradley-Terry Model | 0.22 |

**Discussion**

By fitting and evaluating the Bradley-Terry models, we were able to learn about rankings and approximate the TVI indices for the teams that are used to create the rankings for the NCAA. We, however, found that game location does not have a significant impact on the rankings and the model that approximates the TVI is not able to help determine if Michigan is inferior to other teams. The Bradley-Terry models are specifically designed to be able to look at data that compares two objects in an event and help create a model for rankings. It's commonly used in finding rankings in other sporting events as well, hence it's an obvious choice for this paper as well.

For this study, a big limitation was the data that we analyzed. The data was obtained earlier this year, has very few observations to be able to accurately fit a good model. It would be beneficial for such a study to have access to older data and test the model on the previous rankings for games over a period of years. This would help understand the intricacies of the model and better approximate for the TVI for teams.

Furthermore, it would also be beneficial to look into other methods to account for the game location as a variable. This is important as while our study shows it didn't have a significant impact, it still seems like a practically significant variable to take into account for the

study. Further research on topics and issues addressed in this paper are important to ensure the NCAA rankings are properly and accurately done for the viewers around the globe.

## References

**Hansen, B. (2018).** Paper assignment for Stats Unit 3. *Stats 485: Capstone Seminar*. Retrieved from Course Canvas Site, University of Michigan

**NCAA Viewership and Attendance Numbers** collected from

https://www.ncaa.com/news/basketball-men/article/2018-04-13/2018-ncaa-tournament-and-final-four-viewership-attendance

# Paper3_NCAA_Appendix

*Rifayat Uddin*

*4/12/2019*

## External Requirements

```
### External Requirements
library(ggplot2)
library(MASS)
library(brglm)
```

```
## Loading required package: profileModel
```

```
## 'brglm' will gradually be superseded by 'brglm2' (https://cran.r-project.org/package=brglm2), which p
```

```
library(multcomp)
```

```
## Loading required package: mvtnorm
```

```
## Loading required package: survival
```

```
## Loading required package: TH.data
```

```
##
## Attaching package: 'TH.data'
```

```
## The following object is masked from 'package:MASS':
##
##     geyser
```

```
library(arm)
```

```
## Loading required package: Matrix
```

```
## Loading required package: lme4
```

```
##
## arm (Version 1.10-1, built: 2018-4-12)
```

```
## Working directory is /Users/rifayatuddin/Documents/Stats485/Paper3
```

```
library(tidyverse)
```

```
## -- Attaching packages -------------------------------------------------- tidyverse 1.2.1 --
```

```
## v tibble  2.1.1        v purrr   0.3.2
## v tidyr   0.8.3        v dplyr   0.8.0.1
## v readr   1.3.1        v stringr 1.4.0
## v tibble  2.1.1        v forcats 0.4.0


## -- Conflicts ------------------------------------------------------ tidyverse_conflicts() --
## x tidyr::expand() masks Matrix::expand()
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
## x dplyr::select() masks MASS::select()
```

**library**(caret)

```
## Loading required package: lattice


##
## Attaching package: 'caret'


## The following object is masked from 'package:purrr':
##
##     lift


## The following object is masked from 'package:survival':
##
##     cluster
```

**library**(boot)

```
##
## Attaching package: 'boot'


## The following object is masked from 'package:lattice':
##
##     melanoma


## The following object is masked from 'package:arm':
##
##     logit


## The following object is masked from 'package:survival':
##
##     aml
```

## Overview

This document is the Technical Appendix for Paper 3 Version 1 of Stats 485 Capstone Course. This portion aims at deriving a ranking for the NCAA Men's Basketball Tournament. As the NCAA recently decided to alter the method to evaluate teams by dropping he RPI, this appendix address key questions of ranking, and how rankings change due to fitting with different models and using different coefficients in the evaluations of these models. The appendix shows work in terms of how "game location" could affect rankings and later addresses the "non-inferiority" between Michigan and three teams within the Big Ten to see if Michigan is really inferior to these teams or not.

## Part 1:Bradley Terry Model and Rankings

The data is read and 2 models are fit without the use of RUT in the analysis. This makes calculations easier in the following parts.

```
ncaa_data <- read.csv('https://dept.stat.lsa.umich.edu/~bbh/s485/data/big10menshoopsJan2019.csv')


#Model fitted using glm
bt1 <- glm(home_won ~ . - 1 - RUT, family = binomial, data = ncaa_data)
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
#Model fitted using brglm
bt2 <- brglm(home_won ~ . - 1 - RUT, family = binomial, data = ncaa_data)

#Summary from glm
summary(bt1)
```

```
##
## Call:
## glm(formula = home_won ~ . - 1 - RUT, family = binomial, data = ncaa_data)
##
## Deviance Residuals:
##       Min        1Q    Median        3Q       Max
## -1.17741  -0.00002   0.00000   0.00002   1.17741
##
## Coefficients:
##         Estimate Std. Error z value Pr(>|z|)
## ILL   -6.380e+01  3.647e+04  -0.002    0.999
## IOWA  -2.038e+01  1.619e+04  -0.001    0.999
## IU     5.301e-01  6.841e+04   0.000    1.000
## MARY   2.062e+01  1.825e+04   0.001    0.999
## MINN   8.759e-15  2.000e+00   0.000    1.000
## MSU    2.258e+01  3.437e+04   0.001    0.999
## NEB   -4.179e+01  2.946e+04  -0.001    0.999
## NW    -4.159e+01  2.785e+04  -0.001    0.999
## OSU    2.783e-15  1.732e+00   0.000    1.000
## PSU   -2.227e+01  4.164e+04  -0.001    1.000
## PUR    4.151e+01  2.767e+04   0.002    0.999
## UM     6.267e+01  3.656e+04   0.002    0.999
## WIS    1.542e-14  1.732e+00   0.000    1.000
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 40.2025  on 29  degrees of freedom
## Residual deviance:  5.5452  on 16  degrees of freedom
## AIC: 31.545
##
## Number of Fisher Scoring iterations: 22
```

```
#Summary from brglm
summary(bt2)
```

```
##
## Call:
## brglm(formula = home_won ~ . - 1 - RUT, family = binomial, data = ncaa_data)
##
##
## Coefficients:
##       Estimate Std. Error z value Pr(>|z|)
## ILL  -2.51681    2.13135  -1.181    0.238
## IOWA -0.12242    1.71621  -0.071    0.943
## IU    0.58255    2.24963   0.259    0.796
## MARY  0.79737    1.64744   0.484    0.628
## MINN  0.30972    1.75490   0.176    0.860
## MSU   1.63707    1.74509   0.938    0.348
## NEB  -1.33092    1.91553  -0.695    0.487
## NW   -1.13129    2.03101  -0.557    0.578
## OSU   0.05913    1.49638   0.040    0.968
## PSU  -0.87700    2.02079  -0.434    0.664
## PUR   1.65954    2.07612   0.799    0.424
## UM    2.51307    2.38231   1.055    0.291
## WIS   0.54128    1.49041   0.363    0.716
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 22.673  on 29  degrees of freedom
## Residual deviance: 14.984  on 16  degrees of freedom
## Penalized deviance: 25.19689
## AIC:  40.984
```

```
#Ability coefficients are ranked for each team for each of the models. The higher the ability, the high
rank(coef(bt1))
```

```
##  ILL IOWA   IU MARY MINN  MSU  NEB   NW  OSU  PSU  PUR   UM  WIS
##    1    5    9   10    7   11    2    3    6    4   12   13    8
```

```
rank(coef(bt2))
```

```
##  ILL IOWA   IU MARY MINN  MSU  NEB   NW  OSU  PSU  PUR   UM  WIS
##    1    5    9   10    7   11    2    3    6    4   12   13    8
```

In order to better asses whether the ability coefficients do provide accurate data, a simple test was conducted on the bt1 model to check for accuracy and see if the model correctly classifies team rankings. It was found to be accurate.Rankings remained same for both models, but brglm was used for higher coefficient values found for ease of calculations. The brglm method fits generalized linear models with binomial responses using either an adjusted-score approach to bias reduction or maximum penalized likelihood where penalization by Jeffreys invariant prior. These procedures return estimates with improved frequentist properties (bias, mean squared error) that are always finite even in cases where the maximum likelihood estimates are infinite (data separation), which makes this a better model.

```r
#RANKINGS AND ABILITY TEST
#ILL
ILL_count = 1
#make a vector of the values
vector_coeff = c(-6.379944e+01, -2.038489e+01, 5.301321e-01, 2.062411e+01, 8.759143e-15, 2.258277e+01,

probs_vector = c(0,0,0,0,0,0,0,0,0,0,0,0)
for(i in 2:13) {
  num = vector_coeff[i]
  probs = plogis((-6.379944e+01)-(num))
  #probs
  if(probs > 0.5) {
    ILL_count = ILL_count + 1
  }
}

#IOWA
IOWA_count = 0
vector_coeff = c(-2.038489e+01,-6.379944e+01, 5.301321e-01, 2.062411e+01, 8.759143e-15, 2.258277e+01, -
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    IOWA_count = IOWA_count + 1
  }
}

IU_count = 0
vector_coeff = c(5.301321e-01,-2.038489e+01,-6.379944e+01, 2.062411e+01, 8.759143e-15, 2.258277e+01, -4
#IU
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    IU_count = IU_count + 1
  }
}

MARY_count = 0
vector_coeff = c(2.062411e+01, 5.301321e-01,-2.038489e+01,-6.379944e+01, 8.759143e-15, 2.258277e+01, -4
#IU
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    MARY_count = MARY_count + 1
  }
```

```
}

MINN_count = 0
vector_coeff = c(8.759143e-15,2.062411e+01, 5.301321e-01,-2.038489e+01,-6.379944e+01,  2.258277e+01, -4
#IU
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    MINN_count = MINN_count + 1
  }
}

MSU_count = 0
vector_coeff = c(2.258277e+01, 8.759143e-15,2.062411e+01, 5.301321e-01,-2.038489e+01,-6.379944e+01, -4.1
#MSU
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    MSU_count = MSU_count + 1
  }
}


NEB_count = 0
vector_coeff = c(-4.179402e+01, 2.258277e+01, 8.759143e-15,2.062411e+01, 5.301321e-01,-2.038489e+01,-6.3

for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    NEB_count = NEB_count + 1
  }
}

NW_count = 0
vector_coeff = c( -4.159127e+01, -4.179402e+01, 2.258277e+01, 8.759143e-15,2.062411e+01, 5.301321e-01,-2
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    NW_count = NW_count + 1
  }
```

```r
}

OSU_count = 0
vector_coeff = c(2.782998e-15, -4.159127e+01, -4.179402e+01, 2.258277e+01, 8.759143e-15,2.062411e+01, 5
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    OSU_count = OSU_count + 1
  }
}

PSU_count = 0
vector_coeff = c(-2.227364e+01, 2.782998e-15, -4.159127e+01, -4.179402e+01, 2.258277e+01, 8.759143e-15,
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    PSU_count = PSU_count + 1
  }
}

PUR_count = 0
vector_coeff = c(4.150916e+01,-2.227364e+01, 2.782998e-15, -4.159127e+01, -4.179402e+01, 2.258277e+01,
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    PUR_count = PUR_count + 1
  }
}

UM_count = 0
vector_coeff = c(6.267189e+01,4.150916e+01,-2.227364e+01, 2.782998e-15, -4.159127e+01, -4.179402e+01, 2
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    UM_count = UM_count + 1
  }
}

WIS_count = 0
vector_coeff = c(1.541957e-14, 6.267189e+01,4.150916e+01,-2.227364e+01, 2.782998e-15, -4.159127e+01, -4
```

```r
for(i in 2:13) {
  num = vector_coeff[i]
  int = vector_coeff[1]
  probs = plogis((int)-(num))
  #probs
  if(probs > 0.5) {
    WIS_count = WIS_count + 1
  }
}
```

ILL_count;IOWA_count;IU_count;MARY_count;MINN_count;MSU_count;NEB_count;NW_count;OSU_count;PSU_count;PU

```
## [1] 1
```

```
## [1] 4
```

```
## [1] 8
```

```
## [1] 9
```

```
## [1] 6
```

```
## [1] 10
```

```
## [1] 1
```

```
## [1] 2
```

```
## [1] 5
```

```
## [1] 3
```

```
## [1] 11
```

```
## [1] 12
```

```
## [1] 7
```

```
#UM, PUR, MSU, MARY, IU, WIS, MINN, OSU, IOWA, PSU, NW, NEB, ILL
```

## Part 2: Game-location Advantage

Adapting the model chosen to account for game location to incorporate home-court advantage parameter. An anova-test was done comparing the fit of both the models. P-value not found and most probable cause is it's not significant, hence cannot reject null-hypothesis that both models are essentially similar so adding the game location doesn't change the rankings. bt2 is still chosen since the Residual Deviance is lower and degrees of freedom higher.

```
#bt3 is home-court advantage model
bt3 <- brglm(home_won ~ . - RUT, family = binomial, data = ncaa_data)
summary(bt3)
```

```
##
## Call:
## brglm(formula = home_won ~ . - RUT, family = binomial, data = ncaa_data)
##
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.24274    0.57946   0.419    0.675
## ILL         -1.92122    2.05074  -0.937    0.349
## IOWA         0.03112    1.71797   0.018    0.986
## IU           0.56197    2.17086   0.259    0.796
## MARY         0.84628    1.63396   0.518    0.605
## MINN         0.61228    1.90268   0.322    0.748
## MSU          1.75302    1.74746   1.003    0.316
## NEB         -0.93274    2.00595  -0.465    0.642
## NW          -0.92024    1.95815  -0.470    0.638
## OSU          0.18592    1.50776   0.123    0.902
## PSU         -0.93164    2.02654  -0.460    0.646
## PUR          1.46140    2.02070   0.723    0.470
## UM           2.02051    2.25177   0.897    0.370
## WIS          0.61096    1.53591   0.398    0.691
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 20.647  on 28  degrees of freedom
## Residual deviance: 15.715  on 15  degrees of freedom
## Penalized deviance: 24.01975
## AIC:   43.715
```

```
aov1 = anova(bt2,bt3, test = "LRT")
aov1
```

```
## Analysis of Deviance Table
##
## Model 1: home_won ~ (ILL + IOWA + IU + MARY + MINN + MSU + NEB + NW +
##     OSU + PSU + PUR + RUT + UM + WIS) - 1 - RUT
## Model 2: home_won ~ (ILL + IOWA + IU + MARY + MINN + MSU + NEB + NW +
##     OSU + PSU + PUR + RUT + UM + WIS) - RUT
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1        16     14.984
## 2        15     15.715  1  -0.7306
```

```
rank(coef(bt2))
```

```
##  ILL IOWA   IU MARY MINN  MSU  NEB   NW  OSU  PSU  PUR   UM  WIS
##    1    5    9   10    7   11    2    3    6    4   12   13    8
```

```
rank(coef(bt3))
```

```
## (Intercept)          ILL         IOWA           IU         MARY         MINN
##           7            1            5            8           11           10
##         MSU          NEB           NW          OSU          PSU          PUR
##          13            2            4            6            3           12
##          UM          WIS
##          14            9
```

```
#Tables and Plots will be shown in the paper comparing them
```

## Part 3: Non-Inferiority Tests between Michigan and the following teams; Illinois, OSU and PSU

Conducted non-inferiority tests between Michigan and each of these three teams. More specifically, conduct three simultaneous hypothesis tests of the following form:

```
#avo = brglm(as.factor(home_won) ~ . - 1 - RUT, data = ncaa_data)
#aov_ex = brglm(home_won ~ .- RUT, data = ncaa_data)
num = glht(bt2, linfct = c("UM - ILL <= 0","UM - OSU <= 0", "UM - PSU <= 0"))
pvals = summary(num)$test$pvalues[1:3]
p.adjust(pvals, "holm", n=3)
```

```
## [1] 0.1636073 0.3141239 0.2961316
```

## PROBLEM SET 3

1. Getting rid of Illinois as it's ranked last and removing it won't change the rankings. Hence, this shouldn't affect the appendix values.

```
#Model fitted using glm
bt1 <- glm(as.factor(home_won) ~ . - 1 - RUT - ILL, family = binomial, data = ncaa_data)
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

2. Model 1 with lowest AIC is better, hence prior degrees of freedom = 1 is chosen for the model.

```
mod_1 <- bayesglm(home_won ~ . -1 - RUT - ILL, family = binomial, prior.df= 1, data = ncaa_data)
summary(mod_1)
```

```
##
## Call:
## bayesglm(formula = home_won ~ . - 1 - RUT - ILL, family = binomial,
##     data = ncaa_data, prior.df = 1)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.1922  -0.3086   0.3951   0.6530   1.5138
```

```
## 
## Coefficients:
##      Estimate Std. Error z value Pr(>|z|)
## IOWA  0.50893    1.13172   0.450   0.6529
## IU    1.81526    1.54077   1.178   0.2387
## MARY  1.32630    1.33870   0.991   0.3218
## MINN  1.07920    1.31349   0.822   0.4113
## MSU   3.02030    1.80733   1.671   0.0947 .
## NEB  -0.35777    1.14787  -0.312   0.7553
## NW   -0.09089    1.16586  -0.078   0.9379
## OSU   0.76314    1.11613   0.684   0.4941
## PSU  -1.44251    1.94366  -0.742   0.4580
## PUR   2.46650    1.62573   1.517   0.1292
## UM    4.15544    1.98086   2.098   0.0359 *
## WIS   1.11395    1.19842   0.930   0.3526
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 40.203  on 29  degrees of freedom
## Residual deviance: 16.586  on 17  degrees of freedom
## AIC: 40.586
## 
## Number of Fisher Scoring iterations: 24
```

```r
mod_2 <- bayesglm(home_won ~ . -1 - RUT - ILL, family = binomial, prior.df= Inf , data = ncaa_data)
summary(mod_2)
```

```
## 
## Call:
## bayesglm(formula = home_won ~ . - 1 - RUT - ILL, family = binomial,
##     data = ncaa_data, prior.df = Inf)
## 
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.1832  -0.3094   0.4121   0.7081   1.5610
## 
## Coefficients:
##      Estimate Std. Error z value Pr(>|z|)
## IOWA  0.59228    1.21758   0.486   0.6267
## IU    1.81963    1.56292   1.164   0.2443
## MARY  1.42135    1.42515   0.997   0.3186
## MINN  1.24166    1.40927   0.881   0.3783
## MSU   3.01548    1.77069   1.703   0.0886 .
## NEB  -0.32187    1.22866  -0.262   0.7933
## NW   -0.08688    1.24396  -0.070   0.9443
## OSU   0.86771    1.18476   0.732   0.4639
## PSU  -1.56170    2.15825  -0.724   0.4693
## PUR   2.45903    1.60267   1.534   0.1249
## UM    3.67394    1.67841   2.189   0.0286 *
## WIS   1.25540    1.29601   0.969   0.3327
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 40.203  on 29  degrees of freedom
## Residual deviance: 16.749  on 17  degrees of freedom
## AIC: 40.749
##
## Number of Fisher Scoring iterations: 6
```

```r
mod_3 <- bayesglm(home_won ~ . -1 - RUT - ILL, family = binomial, prior.df= 1000 , data = ncaa_data)
summary(mod_3)
```

```
##
## Call:
## bayesglm(formula = home_won ~ . - 1 - RUT - ILL, family = binomial,
##     data = ncaa_data, prior.df = 1000)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.1833  -0.3094   0.4121   0.7082   1.5609
##
## Coefficients:
##       Estimate Std. Error z value Pr(>|z|)
## IOWA   0.59219    1.21747   0.486   0.6267
## IU     1.81969    1.56292   1.164   0.2443
## MARY   1.42125    1.42504   0.997   0.3186
## MINN   1.24145    1.40915   0.881   0.3783
## MSU    3.01552    1.77075   1.703   0.0886 .
## NEB   -0.32192    1.22855  -0.262   0.7933
## NW    -0.08686    1.24387  -0.070   0.9443
## OSU    0.86758    1.18467   0.732   0.4640
## PSU   -1.56152    2.15796  -0.724   0.4693
## PUR    2.45911    1.60272   1.534   0.1249
## UM     3.67476    1.67884   2.189   0.0286 *
## WIS    1.25522    1.29587   0.969   0.3327
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 40.203  on 29  degrees of freedom
## Residual deviance: 16.749  on 17  degrees of freedom
## AIC: 40.749
##
## Number of Fisher Scoring iterations: 7
```

3. Finding K; K = 10 gives a pretty good RMSE, which is what was chosen. This also prevents the problem posed in the question.

```r
# Define training control
set.seed(123)
train.control <- trainControl(method = "cv", number = 10)
# Train the model
```

```
model <- train(as.factor(home_won) ~ . -1 - RUT - ILL, data = ncaa_data,  method = "bayesglm", trControl
# Summarize the results
print(model)
```

```
## Bayesian Generalized Linear Model
##
## 29 samples
## 14 predictors
##  2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 27, 26, 27, 26, 26, 26, ...
## Resampling results:
##
##   Accuracy  Kappa
##   0.65      0.23
```

```
# Define training control
set.seed(123)
train.control <- trainControl(method = "cv", number = 5)
# Train the model
model <- train(as.factor(home_won) ~ . -1 - RUT - ILL, data = ncaa_data,  method = "bayesglm", trControl
# Summarize the results
print(model)
```

```
## Bayesian Generalized Linear Model
##
## 29 samples
## 14 predictors
##  2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 24, 23, 24, 22, 23
## Resampling results:
##
##   Accuracy   Kappa
##   0.6961905  0.3803463
```

Having found K = 10, we compare the models (mod_1 is for the bayesglm and bt1 using glm)

```
cv_error = cv.glm(ncaa_data, mod_1, K = 10)
cv_error$delta[1]
```

```
## [1] 0.2305938
```

```
cv_error = cv.glm(ncaa_data, bt1, K = 10)
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```r
cv_error$delta[1]
```

```
## [1] 0.25922
```

```r
cv_error = cv.glm(ncaa_data, bt2, K = 10)
cv_error$delta[1]
```

```
## [1] 0.2059884
```

```r
rank(coef(bt2))
```

```
##  ILL IOWA   IU MARY MINN  MSU  NEB   NW  OSU  PSU  PUR   UM  WIS
##    1    5    9   10    7   11    2    3    6    4   12   13    8
```

The error for Bayesglm is lower, hence we decide to use that model.