

CATEGORICAL ANALYSIS OF LOYOLA RAMBLERS GAMES (2016 - PRESENT), VERSION 2.0

CHARLES HWANG (LUC '20)

BACKGROUND/MOTIVATIONS

- **UNIQUENESS**
 - ONLY ONE PREVIOUS KNOWN STATISTICAL ANALYSIS ON THESE DATA HAD BEEN PERFORMED
- **PERSONAL INTEREST**
- **EXISTING DATASET**
- **MORE TIME FOR QUESTIONS**

CRITERIA FOR DATASET INCLUSION

- ALL MATCHES BETWEEN AUGUST 19, 2016 AND SEPTEMBER 12, 2022
- LOYOLA RAMBLERS' PROGRAMS PLAYING "HEAD-TO-HEAD" COMPETITION INCLUDE WOMEN'S SOFTBALL (SB) AND BOTH WOMEN'S AND MEN'S SOCCER (WSOC, MSOC), VOLLEYBALL (WVB, MVB), AND BASKETBALL (WBB, MBB)
 - WOMEN'S AND MEN'S CROSS COUNTRY, TRACK & FIELD, GOLF, AND SPIRIT TEAMS EXCLUDED

VARIABLES INCLUDED

- TEAM
- SEASON
- TYPE (REGULAR SEASON, CONFERENCE TOURNAMENT, POSTSEASON, ETC.)
- LOCATION (HOME, AWAY, OR NEUTRAL)
- OPPONENT
- OPPONENT CONFERENCE
- OUTCOME (WIN, DRAW, OR LOSS)
- DAY OF WEEK
- MONTH
- SCORE FOR
- SCORE AGAINST
- LENGTH (REGULATION VS. OVERTIME, ETC.)
- OUTCOME OF PREVIOUS GAME
- WIN/LOSS STREAK PRIOR TO GAME
- WHETHER OR NOT I ATTENDED

EXPLORATORY DATA ANALYSIS (EDA)

- ALL DATA WERE COMPILED MANUALLY IN EXCEL THEN READ INTO R AS A .CSV FILE
 - N = 1,203 GAMES
- VARIABLES SELECTED FROM THOSE EASILY AVAILABLE ONLINE¹
- CONFERENCE AT TIME OF GAME RECORDED^{2, 3}
- DIFFERENT CONFERENCES PER SPORT²
- ACCURACY CHECKS
 - MISSING DATA (NONE)
 - IMPOSSIBLE VALUES
 - SIMILAR VALUES (NORTHERN IOWA VS. NORTHERN ILLINOIS, MVC VS. MAC, ETC.)

HYPOTHESES/RESEARCH QUESTIONS

1. WHICH VARIABLE(S) ARE THE MOST SIGNIFICANT FOR EACH TEAM?
2. DID EACH TEAM PERFORM DIFFERENTLY DURING THE SEASONS BEFORE, AFTER, AND THE SUCCESSIVE SEASON AFTER THE ONSET OF THE COVID-19 PANDEMIC?
3. DID TEAMS PERFORM DIFFERENTLY BEFORE AND AFTER A HEAD COACHING CHANGE?
 - WVB (AFTER 2017), SB (AFTER 2019), MBB (AFTER 2020-21)
4. DO TEAMS PERFORM DIFFERENTLY AT HOME THAN AWAY/NEUTRAL?
5. DO TEAMS PERFORM DIFFERENTLY WHEN I ATTEND?

PREDICTIONS

1. LOCATION IS SIGNIFICANT
 - SEASON AND WIN/LOSS STREAK MAY BE SIGNIFICANT
2. TEAMS GENERALLY PERFORMED DIFFERENTLY IN EACH OF THE THREE SEASONS
3. WVB DID BETTER, SB AND MBB HAD NO STATISTICALLY SIGNIFICANT DIFFERENCE
4. TEAMS PERFORM BETTER AT HOME THAN AWAY
5. TEAMS PERFORM BETTER WHEN I ATTEND

ANALYSES PERFORMED

1. GLM (STEPWISE SELECTION WITH AIC, BOTH DIRECTIONS)

2-3. (A LOT OF) TWO-PROPORTION Z-TESTS ($\alpha = 0.05$)

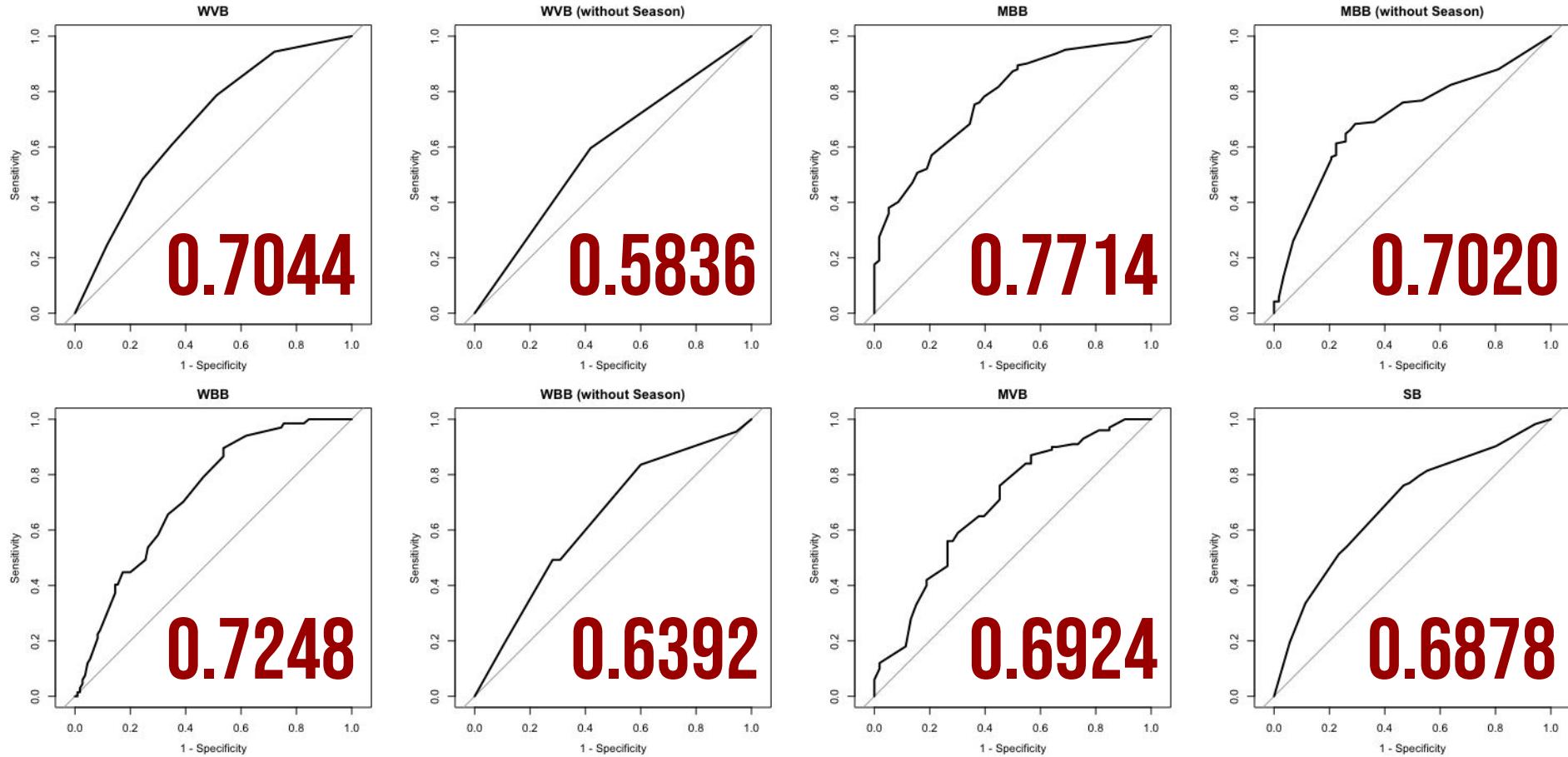
4-5. χ^2 -TESTS OF INDEPENDENCE ($\alpha = 0.01$)

5. LIKELIHOOD RATIO χ^2 G-TEST, FISHER'S EXACT TEST

(NO MSOC OR WSOC)	SEASON	TYPE	LOCATION	MONTH	MATCH LENGTH	PREVIOUS OUTCOME	WIN/LOSS STREAK #	DID I ATTEND?	STEPS
WVB	'17, (2)		<u>HOME, N</u>						8
WVB*							<u>0.023603</u>		8
WBB	(5)						<u>0.080987</u>		7
WBB*			<u>HOME, N</u>			<u>0.008321</u>			6
MBB	(3)		<u>HOME, N</u>					<u>0.084409</u>	8
MBB*			<u>HOME, N</u>				<u>0.121238</u>		8
MVB	(0)	REGULAR			<u>4, 5 SETS</u>			<u>0.055525</u>	5
SOFTBALL			<u>HOME, N</u>	<u>0.017802</u>					9

*WITHOUT SEASON VARIABLE

SCORE WAS EXCLUDED FROM ALL, OPPONENT AND CONFERENCE WERE EXCLUDED FROM SOME



2. TWO-PROPORTION Z-TESTS ON WIN % (YEARS 0, 1, AND 2 OF PANDEMIC)

	WSOC	MSOC	WVB	WBB*	MBB	MVB*	SB*
0 VS. 1	0.77139	0.72546	0.44357	0.58592	<u>0.09625</u>	<u>0.04106</u>	<u>0.04518</u>
0 VS. 2	0.96522	0.82825	0.92523	0.52211	0.36922	<u>0.06620</u>	0.10618
1 VS. 2	0.78945	0.87174	0.37403	0.24023	0.42010	0.72402	0.60318

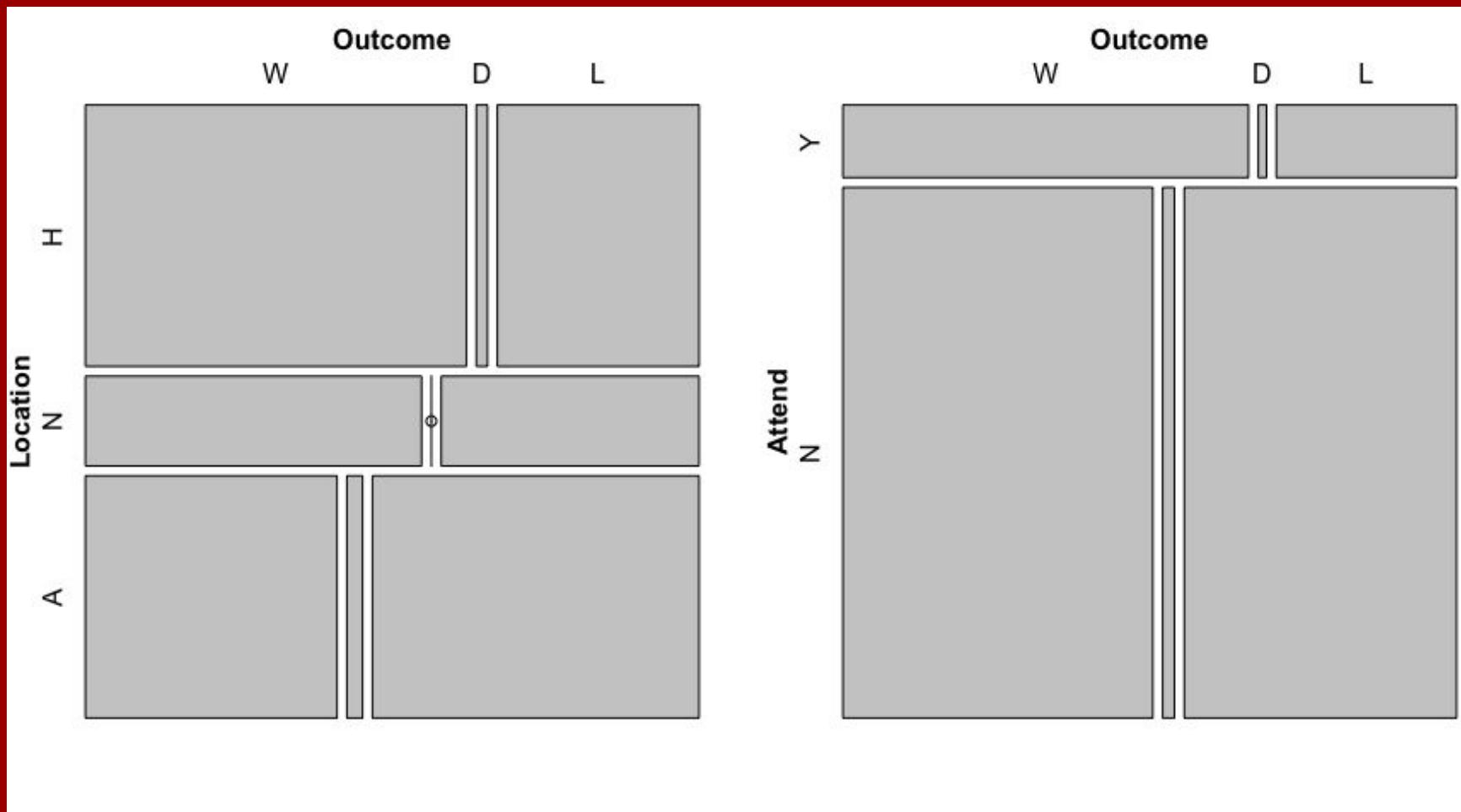
*1+ WBB, 9+ MVB, 31+ SB GAMES CANCELLED IN 2020

3. TWO-PROPORTION Z-TESTS ON WIN % (BEFORE & AFTER HC CHANGE [†])

	X-STATISTIC*	P-VALUE	WIN % BEFORE	WIN % AFTER	95% C.I.
WVB (2017)	14.87044	<u>0.0001151535</u>	0.3114754	0.6129032	{-0.446,-0.157}
SB (2019)	3.043565	<u>0.08105772</u>	0.4901961	0.3818182	{-0.013, 0.229}
* $\chi^2(1) = Z^2$					
[†] NBB (2021) 64340527 GAMES FOR MEN, 0.4398121 WOMENS BASKETBALL, 0.7005688 GAMES FOR MEN				0.7575758	{-0.219, 0.105}
VOLLEYBALL					

4-5. TWO-WAY FREQUENCY TABLES

	HOME	NEUTRAL 1	AWAY		YES	NO	TOTAL
WIN	340	103	208	WIN	99	552	651
DRAW	10	0	13	DRAW ²	2	21	23
LOSS	180	79	270	LOSS	44	485	529
TOTAL	530	182	491	TOTAL	145	1,058	1,203
WIN %	0.65094	0.56593	0.43686	WIN %	0.6897	0.5317	0.5507



4-5. TWO-WAY FREQUENCY TABLES (MODIFIED)

	HOME	N+AWAY ¹
WIN	340	311
DRAW	10	13
LOSS	180	349
TOTAL	530	673
WIN %	0.65094	0.47177

	YES	NO	TOTAL
WIN ²	99	573	672
LOSS ²	46	485	531
TOTAL	145	1,058	1,203
WIN %	0.6828	0.5416	0.5586

ODDS RATIO: $(99*485)/(573*46) \approx 1.82164808$

4-5. χ^2 -TESTS OF INDEPENDENCE, LR χ^2 G-TEST, FISHER'S EXACT TEST

	OUTCOME VS. LOCATION	OUTCOME VS. LOCATION ¹	LOCATION (LR χ^2 G-TEST ¹)	OUTCOME VS. ATTENDANCE	OUTCOME VS. ATTENDANCE ²	ATTENDANCE ² (FISHER'S)
STATISTIC	53.11884	39.22969	39.59474	13.31843	10.3073	
P-VALUE	< 0.0000001	0.00000000 6	0.00000000 3	0.00128215 2	0.00132505 3	0.00129913 8
SIGNIF. CODE	***	***	***	**	**	**

¹NEUTRAL AND AWAY GAMES GROUPED TOGETHER

²ATTENDED DRAWS COUNTED AS LOSSES, NON-ATTENDED DRAWS COUNTED AS WINS

CONCLUSIONS

1. LOCATION IS THE MOST SIGNIFICANT VARIABLE OVERALL (SEE #4) ✓
 - SEASON IS ALSO SIGNIFICANT
2. MOST TEAMS DID NOT PERFORM DIFFERENTLY BEFORE AND AFTER THE PANDEMIC X
 - EXCEPTIONS: MVB ($P = 0.04106$) AND SB ($P = 0.04518$), 2020 VS. 2021
3. WB DID BETTER, MBB HAD NO DIFFERENCE, SB DID WORSE (NOT SIGNIFICANT) ✓
4. TEAMS DID BETTER AT HOME ($P < 0.00000001$) ✓
5. TEAMS DID BETTER WHEN I ATTENDED ($P = 0.001310144$) ✓

ANALYSIS

1. SEE #4, PERFORMANCE BETWEEN SEASONS BEING DIFFERENT ALSO MAKES SENSE
 - OTHER VARIABLES (PREVIOUS MATCH, MATCH LENGTH, MONTH) BEING STATISTICALLY SIGNIFICANT FOR INDIVIDUAL PROGRAMS IS INTERESTING
2. THIS IS INTERESTING BECAUSE I EXPECTED TRAINING RESTRICTIONS, STRENGTH & CONDITIONING, MATCH CONDITIONS (E.G., NO FANS), ETC. TO BE DIFFERENT, BUT IT MAKES SENSE THAT THEY ARE NOT STATISTICALLY SIGNIFICANT AT $\alpha = 0.05$
3. THIS WAS EXPECTED AS I FOLLOWED ALL THE TEAMS BEFORE AND AFTER
 - SOFTBALL GETTING “WORSE” COULD BE ATTRIBUTED TO NON-CON SCHEDULING
4. THIS MAKES SENSE, MOST TEAMS PERFORM BETTER AT HOME

ANALYSIS

5.

RANDOMNESS AND INDEPENDENCE ASSUMPTIONS

- MY ATTENDANCE OF MATCHES WAS NOT RANDOM
 - INTUITIVELY, FANS ATTEND MATCHES IF THEY WANT TO
 - LIKELIHOOD OF ATTENDANCE MAY BE DEPENDENT ON PAST ATTENDANCE, GAMEDAY PROMOTIONS/GIVEAWAYS, DISTANCE, ETC.
- TEAM PERFORMANCE IS LIKELY NOT INDEPENDENT
 - PERFORMANCE IN A GIVEN MATCH MAY DEPEND ON PREVIOUS MATCHES
- HOWEVER, IT APPEARS THESE ASSUMPTIONS ARE NOT “BADLY” VIOLATED.

IMPROVEMENTS MADE

- FURTHER ANALYSIS ON THE UNDERLYING REASONS FOR RESULTS ✓
- REVIEWING TEST ASSUMPTIONS (RANDOMNESS, INDEPENDENCE, ETC.) ✓
- UPDATING DATASET AND RESULTS WITH NEW DATA

FUTURE IMPROVEMENTS

- ADDING OTHER VARIABLES AND INTERACTION TERMS
 - PRIVATE VS. PUBLIC, DAYS SINCE PREVIOUS GAME, DATE, TIME, DURATION, “OFFICIAL” ATTENDANCE, NATIONAL POLL RANKINGS, DISTANCE FROM CAMPUS, DISTANCE TRAVELED, TICKET PRICE, WHETHER THERE WAS A PROMOTION/GIVEAWAY, ETC.
- TESTS ON NON-CONFERENCE AND CONFERENCE RECORDS INSTEAD OF OVERALL RECORD
- ADDING NON-HEAD-TO-HEAD TEAMS (WITH FINAL RANKING/POSITION AS OUTCOME)
- RESEARCHING HOW TO CODE DRAWS IN GLM FOR SOCCER
- PREDICTIVE ANALYTICS (CART, RANDOM FOREST, ETC.), MACHINE LEARNING
- TIME-SERIES ANALYSIS

SPRING 2022 PROJECT REPORT (3 PAGES) AND CODE APPENDIX (19 PAGES)

Analysis of Loyola Ramblers Games (2016 - present)

Charles Hwang

4/25/2022

The dataset I chose for this project is the game data for all head-to-head Loyola Ramblers programs from the beginning of the Fall 2016 through April 24, 2022. This includes women's softball and both women's and men's soccer, basketball, and volleyball, but excludes competitions with multiple teams by nature (cross country, track and field, golf, darts, chess, etc.). There were 1,171 games from the various public sources that I gathered from the athletic department websites (<https://www.loyolahawks.com>). There are 1,171 games in total. I recorded 15 variables: (1) team, (2) season, (3) type of game (regular season, conference tournament, postseason, etc.), (4) location (home, away, or neutral), (5) opponent, (6) opponent conference, (7) day of week, (8) month, (9) outcome (win, draw, or loss), (10) score for and (11) against, (12) length (regulation vs. overtime, etc.)¹, (13) outcome of previous match, (14) win/loss streak prior to game², and (15) whether or not I attended.

If (5) opponents' names changed, they were entered under the same name if they were functionally representing the same university.³ One of the many challenges faced was regarding (6) conference, since conferences vary by sport (meaning athletic departments are in different conferences depending on the sport), especially in men's volleyball and men's soccer. Several opponents also changed conferences between the game and this project, so I recorded the opponent's conference at the time of the game. Non-Division I opponents were recorded as the division/association and the conference (for example, an NCAA Division II opponent in the GLVC is "DIIGLVC" and a NAIA opponent in the CAC is "NAIACCAC"). Finally, some opponents were independent (their team did not have a conference) and I recorded them as such. (9) Month was coded as a number for potential quantitative analysis if desired.

After manually entering variables 1-12 and 14-15 in an Excel spreadsheet, I conducted some basic accuracy checks, like ensuring there was no missing data and that (3) type, (4) location, (8) day, (10) outcome, and (15) attendance had all the possible levels listed. I also checked to make sure the opponents and conferences were correctly by sorting and inspecting the columns. After exporting the spreadsheet as a .csv file and reading it into R, I set several variables as factors. I also split the (11) streak variable and created a new variable for (13) previous outcome (Figure 1).

Going into the project, I didn't have any specific questions or hypotheses to explore; I was mainly planning on creating and analyzing the dataset and finding interesting things to look at. After doing so, the main questions I had were:

- Which variable(s) are the most significant for each team?
- Did each team perform differently during the seasons before, after, and the successive season after the onset of the COVID-19 pandemic?

¹This varies by sport. For soccer, NCAA rules (prior to this month) stipulated two 10-minute periods of sudden death overtime (and penalty kicks afterwards for conference tournament and postseason matches), so the levels were "FT" (regulation), "OT1", "OT2", and "PK". Volleyball matches are best-of-five sets, so I entered the number of sets (3, 4, or 5). Basketball uses five-minute halves, so the levels were "OT", "OT1", "OT2", and "OT3", etc. Softball has a mercy rule after 5 innings, so games can actually be shorter or longer than 7 innnges.

²I used streak prior to game because I did not feel the streak after the game affected the game itself. It was the first game of the season. I also included streak from the previous season.

³For example, IPFW became Fort Wayne in 2016 and Purdue Fort Wayne in 2018, and all instances of them in the data were recorded as "Purdue Fort Wayne".

- For teams that had a change in the head coaching position during the timeframe of the data, did those teams perform differently before and after the?
- Do teams perform differently at home than away?
- Do teams perform differently when I attended?

I first decided to create generalized linear models (GLM) using stepwise regression for each team and (9) outcome as the dependent variable to see which variables are the most significant. I ran models for each team except soccer (more information in the conclusion) and the results can be found in Figures 2 through 6. I excluded (10) score for and (11) against because they would be perfect predictors of (9) outcome. I also excluded (5) opponent and (6) conference from several models because there were too many levels in each to be practical in a model. I found that several levels for (2) season were in almost all of the final models, which suggest season is a factor in determining performance. I ran a logistic regression excluding (2) season from the initial model and found the final models to be different (Figures 2b, 3b, and 4b). Based on prior knowledge of each team, most of the variables in the models made sense to me, but further analysis and discussion of the reasoning for each model would be too lengthy to include in this paper.

I then checked to see whether teams performed differently before, after, and the season after the onset of the COVID-19 pandemic.⁴ Since win percentage is between 0 and 1 and there were three different percentages for each team, a "round-robin" of three two-proportion z-test for each team appeared to be the most effective test. After running all $7 \times 3 = 21$ tests, I found that most teams did not perform differently between these three seasons (Figures 7 through 13). The exceptions are with the men's volleyball ($p = 0.04106$) and softball ($p = 0.040518$) teams, both between the 2020 and 2021 seasons. However, the 2020 seasons of both teams were suspended due to the pandemic which cancelled games that may have changed this result. The 2020 men's volleyball team also performed poorly prior to the suspension which may not be attributable to the pandemic.

In answering the third question, three different teams⁵ had a change in the head coaching position during the timeframe of this data: women's volleyball after the 2017 season, softball after the 2019 season, and men's basketball after the 2021 season. I again used two-proportion z-tests to compare the win percentages of all games under the new coach to the win percentages of all games under the old coach. I found that for the women's volleyball, the win percentage for the men's volleyball increased (Figure 14a), the win percentage for the men's basketball did not have any statistically significant change (Figure 14c), and the win percentage for the softball team actually decreased (Figure 14b), but not at a statistically significant level ($p = 0.07368$). Since the 2020 season was the first season under the new head coach for softball, this may have served as a proxy for pre- and post-pandemic performance.

I seemed apprised with the fourth question that teams perform better at home than away, but I wanted to see the statistical strength of this relationship/association. Since there were two variables, I ran a χ^2 -test on (9) outcome and (4) location with the entire dataset and found that there was strong evidence teams performed better at home ($p < 0.000000001$) and I created a mosaic plot to illustrate this (Figure 16). However, R produced a warning for the `chisq.test()` function and I noticed in the two-way frequency table (Figure 15) that there were $0 < 5$ draws at a neutral venue, which violated the expected value assumption for a χ^2 -test. In thinking about how to proceed, I considered several options in reordering the data but ultimately chose to count neutral games as away games, as they tend to act like "away" games for both teams in terms of facilities, travel, familiarity, fan support, etc. I ran the χ^2 -test (Figure 17) and found similar strong evidence of better team performance at home ($p < 0.000000001$).

Finally, the final and potentially most controversial hypothesis question was posed as personal curiosity but also a form of validation of some sorts. The (9) outcome and (15) attendance variables were again presented most practically as a frequency table and the χ^2 -test (Figure 18) confirmed my intuition about helping teams perform better ($p = 0.0007813$). The mosaic plot (Figure 19) also visualized this significant difference. However, the expected value assumption was again violated as I have attended $0 < 5$ draws in my six years at Loyola.

⁴For soccer and women's volleyball, this was the 2019, 2020, and 2021 seasons. For basketball, this was the 2019-20, 2020-21, and 2021-22 seasons. For men's volleyball and softball, this was the 2020, 2021, and 2022 seasons.

⁵There were also head coaching changes in men's soccer and women's basketball in 2022, but neither team have played any games under the new head coach yet.

Loyola (as I have not attended many soccer games). The nature of the test meant I had to reallocate draws and I eventually decided to list them as wins when running the modified test (Figure 20), which output the same result ($p = 0.001271$). Because of the complications with assumptions, I also ran a Fisher's exact test on the same data (Figure 21) which concurred with the previous χ^2 -tests ($p = 0.000131$).

Some of the advantages that I realized when compiling this dataset was a familiarity and interest in the data. I was knowledgeable about the subject which made it easier to determine variables and analyses. Almost all of the data were also all publicly and conveniently available with a minimal number of clicks which made the process easier. Other athletic departments do not maintain accurate historical records or statistics of games from previous seasons, due to lack of resources or indifference, which can make it difficult to find information from older games.

One of the main disadvantages of this process included the manual labor involved in inputting all the data, which was tedious and took several hours but also risked human error in incorrectly recording variables and getting distractred, etc. Additionally, the athletic department website has been removed a few times since 2016, so the data may not be as reliable as initially intended. For example, in checking for (13) length, summaries of some games will show the score and whether there was overtime, but others (especially softball) are not listed in the score and instead only shown when the row for the game was expanded. Although the data were still publicly available, it took extra time to expand the row to check for length which can compound when checking so many games, and I did not discover this until midway through which could have led to some inaccurate data. The links to recaps and box scores prior to 2017 also tended to be dead returning a 404 error.

Other disadvantages included not remembering whether I attended a specific game. Although I previously remembered specific details of every game I attended (I could recall the score, statistics, emotions, etc.) prior to the pandemic, I never recorded this information anywhere and two years had passed since this information was useful/practical to me. I know the number of games I attended⁶ but struggled to recall whether I attended certain games and had to verify video, class schedules, social media, etc. which took a lot of extra time. Some of the opponent logos for older games (and rarely, some statistics) were also incorrect which led to some initial confusion.

Future analysis could include additional variables, like days since previous game, date, time, duration, attendance, national poll rankings, distance from Loyola, ticket price, whether there was a promotion/giveaway, etc. Additional statistical tests could also be run. I do not believe I was able to find a way to code draws for soccer which is why there were no GLM for those teams. Additional research or different statistical methods could produce suitable models for them. I also realized during the project that the randomness and independence assumptions for some of the tests may not be completely satisfied. Lastly, I felt I didn't have enough time to explain the reasoning behind the results, although I could infer it based on context and prior knowledge when reviewing the output.

I felt I chose several different types of variables (factors with different numbers of levels and quantitative) to record and analyze. Despite the high amount of work spent during a short time period, this project was a good categorical analysis of the data.

Appendix

Figure 1

```
rs(list=ls)
L<-read.csv("Users/neuser/Desktop/Notes/Graduate/STAT 410 - Categorical Data Analysis/LUC.csv",header=TRUE)
L$Sport<-as.factor(L$Sport)
L$Season<-as.factor(L$Season)
L$Type<-as.factor(L$Type)
L$Location<-as.factor(L$Location)
```

⁶Officially 136, because I chose to only count games if I was present at the end of the game (i.e., not counting games if I left early), with some discretion. However, the final count in the spreadsheet is 137.

REFERENCES

1. LOYOLARAMBLERS.COM (AND ITS THOUSANDS OF SUBPAGES)
2. OTHER TEAMS' ATHLETICS WEBSITES (TO VERIFY OLDER MATCHES AND CONFERENCE AFFILIATIONS)
3. WIKIPEDIA (TO VERIFY CONFERENCE AFFILIATIONS)

QUESTIONS?

YOU JUST GOT

RAMBLED