Prediction and Prevention Metrics and Statistics Aimed at Understanding UCL Injury Risk

Will Rettig

11/23/2025

Denison University

Table of Contents

Abstract

The game of baseball, forever known as America's pastime, is enveloped in a pandemic. Arm

injuries continuously plague athletes of all ages on the diamond. Elbow injuries cover roughly

20% of all injuries in professional baseball (Mine et al., 2021). Statistics on the frequencies of

UCL injuries have skyrocketed over the past 10 years and more (Jang, 2019). While there are

surgeries and rehabilitation procedures in place for recovery, there is little currently being done

to consider the prevention and prediction of such injuries. Baseball analytics have exploded in

recent years and there is now more readily available data on individuals. This data could play

crucial in determining potential prediction of future injury as well as exposing patterns and

trends that may have led to previous injury. The analysis of this paper will explore metrics,

trends, and insights seeking to uncover commonalities between baseball players with the purpose

in mind of reducing the number of UCL injuries requiring Tommy John surgery.

Prediction and Prevention Metrics and Statistics Aimed at Understanding UCL Injury Risk

**Introduction**

**Background**

The history of arm injuries in the game of baseball can best be classified as dynamic. Arm injuries have varied by location, severity, and prevalence as the years go on and as ideologies change within the game. Trends in today's game show a growing issue, a pandemic, of UCL injuries. The ulnar collateral ligament, known as the UCL, is a sacred ligament in the game of baseball. Located on the inside of the elbow, it holds several tissues together that are important to the throwing motion. Most commonly, baseball players will tear or stretch this ligament due to overuse, or simply greater stress on the ligament (*Ulnar Collateral Ligament (UCL) Injury*, 2022). It just so happens that the new age of baseball has begun, commonly referred to as the 'Statcast Era,' in Major League Baseball. Statcast was first introduced in 2015, with expanded statistics and metrics on batted and pitched ball data, including information collected on every single pitch and swing within a game (Major League Baseball, n.d.). The following trend has been the result; with more numbers available, there is a greater ability to understand what must be done to throw harder. Hitters and pitchers constantly look for an advantage over the other. For hitters, the stronger they can be, the further they can hit the ball. The better their strength and the higher their launch angle, the more likely they are to hit home runs. Inversely, the way that a pitcher must combat this shift is to make the hitter miss their pitches more often. This is best done by two options, the first being creating more spin and being more 'nasty,' and the second being able, simply, to throw harder. The lengths taken to increase spin generation and pitching

velocity have adverse effects on arm health. The body and the arm must move faster and exert more force, which is inherently adding more stress to the arm, specifically the UCL. Further, pitchers are willing to make the sacrifice of putting their arm health at risk if it may lead to advantages over a hitter.

In this paper, data will be explored regarding arm injuries with the purpose of underpinning deficiencies and commonalities that may lead to injury prediction and, ultimately, injury prevention. The research question is best defined as, "can arm injuries, specifically injuries to the UCL, be prevented or predicted through analysis of various data points and systems that expose indicators of weakness or potential injuries?" Even in a sport where arm injuries vary by type and severity, this remains an ever-important question because baseball players of all ages are suffering UCL injuries (Petersen, 2025). These injuries have been known to consistently disable an athlete from returning to their full potential, often cutting a player's career short and can even force early retirement altogether. However, the surgery process has become so successful regarding repair that the injury has become nearly obsolete in terms of the athlete's ability to get back to their full potential. Nevertheless, an athlete loses 12-16 months, on average, of baseball progress with rehabilitation and recovery from the surgery process. This does not take into consideration the mental toll it takes on the injured player. Trends of deficiencies and commonalities would be crucial to reducing the prevalence of this major injury. It would be significant to find alternative methods and approaches to arm care, rehabilitation, and injury prevention, which would outline major implications for the game of baseball. Countless coaches and instructors offer methods and training strategies to 'bullet-proof the UCL,' yet very little science and numbers exist to support these routines. From this excess of data, it is natural to hypothesize that there is some way to predict arm injuries or areas of weakness based on current

data available. Finding the specific location of these areas of weakness are highlighted will only be exposed through analysis and interpretation of the methods and data at hand.

      The needs dictated by the research question, 'can arm injuries be prevented or predicted through analysis of various data points and systems that expose indicators of weakness or potential injuries,' do not explicitly explain a cause-and-effect relationship. Because each athlete has a different genetic makeup, arsenal of pitches, and mechanical foundation and breakdown, the results will be difficult to generalize to an individual. This lack of specific generalization does not indicate an insufficient research question, rather it is important for exposing the trends on a macro-scale as opposed to an individual or micro-scale. Both micro and macro data are utilized for the analysis and within methods, yet generalizations will remain on a wider, broad scale. Further, when individual analysis of an athlete is available, it is important to examine trends, metrics, and statistics as they pertain to the specific individual. When discussed in explicit terms, generalizability will be exclusive to the overall trends for the specific athletes selected by Alpha Baseball and Major League Baseball pitchers in the Statcast Era (2015-present). The results will not be sufficient in proving causality due to the individual differences between players. There are no methods of which this process can enhance and enable causality. This analysis is aimed to see if overlap exists between different data and player sources for variables between platforms and whether these trends can be connected in any sort of way. The exploratory element along with the fundamental analysis of the connection is necessary to future development and research. This analysis is the first of its kind to establish the connection of these various data platforms, further emphasizing the importance of this research and interpretation.

**Literature Review**

While the current discussed analysis is the first of its kind, it is important to note that significant research and analyses exist on the previously referenced history of arm injuries. Though different, these other sources remain important to highlight for the bounds of this analysis. The domain knowledge report attempts to encapsulate and provide an overview into the past and present relevant scholarly information on the topic of arm injuries in baseball players. Much is known regarding the severity of these injuries, yet researchers are skeptical on the true factors that lead into major injuries, such as tears (partial or full) of the ulnar collateral ligament, tears of the rotator cuff, etc. One example of a commonly held theory is that the extreme wear and tear on the arm that causes an injury to the UCL is exacerbated by deficiencies in other areas. These potential deficiencies will be explored using the micro data obtained from Alpha Baseball. These deficiencies are, but not limited to, insufficient hip mobility, insufficient shoulder range of motion, insufficient internal and/or external rotation of the elbow, etc. Another area of focus for instruction is to increase weight and body mass, regardless of player height. Each of these factors are commonly treated with equal importance, which furthers the importance of the research behind them to pinpoint and understand potential causes of significant arm injuries. However, the lacking inclusion of a holistic approach with all potential factors into one major, relevant study is alarming.

Gaining sufficient knowledge of the relevant discussion in this field is important to gaining a true understanding of how the discussed holistic approach can be implemented. The following domain knowledge studies will be listed in chronological order with emphasis on showing the growth and dynamicity of research in the field. (Brown et al., 1988) examined the range of motion of the upper extremities including the shoulder and elbow. The measurements

were taken from MLB baseball players. This study was conducted utilizing measurement with an

emphasis on degrees for internal and external rotation and time used paired with the degrees

allowed for the calculation of force. There were both position players and pitchers utilized in this

study to find the differences between the two groups. It is an important implementation having

both position and pitching players included. A small section of the exploratory analysis of this

study includes position players, yet there are fundamentally less arm metrics and statistics

gathered from position players due to the lack of UCL injury prevalence in position players.

Internal and external elbow rotation deficiencies were not found to be significant as they pertain

to UCL injuries, nor is data readily available on MLB pitchers and position players on these

specific elbow metrics. (Bartlett et al., 1989) examined the force production of the upper

extremities and its parts along with the correlation to the throwing speed of overhead athletes.

This study seems inherently flawed, as it seems logical to conclude without any evidence that

higher force would naturally lead to greater throwing velocities, as measured by a radar gun.

Another less commonly considered flaw is the radar gun itself, as even in 2025 there is only one

pitch speed measuring system known to be the most accurate. Depending on the angle, distance

from the pitcher, and other factors, the speeds can vary by three to five MPH. Pitch speeds, then,

are subjective depending on the radar gun utilized. Regardless of the flaws, this study adds the

following notion to the domain knowledge; more force production does equal higher throwing

velocities. Pitch speed metrics were collected with stadium Trackman, accessible by all 30 MLB

teams, and the top fastball velocities from the Alpha Baseball study were collected using the

same Trackman system. These results were examined also in the exploratory analysis phase,

testing if higher force could result in higher wear and tear on the upper extremities, which is

enticing to look at the relationship between higher velocities and injury prevalence in high

performance athletes. (Wilk et al., 2002) examined the rehabilitation techniques for overhead

throwing athletes at the time. Included are significant analyses, paired with an overview of

potential rehabilitation techniques. There are also areas of potential deficiencies that could lead

to a need for rehabilitation. This scholarly discussion is an important, as it lists and identifies all

the factors that should be considered for a holistic analysis. It provides a sound basis for a

holistic approach to potential research and further emphasizes the need to examine factors on a

macro-level. (Scher et al., 2010) emphasized the importance for examination of upper half and

lower half measurements. It attempts to examine relationships between hip range of motion and

shoulder range of motion and the potential impacts of deficits on injury history. This is an

important discussion as it studies the relationship of both upper half and lower half

measurements, and how the relationship can impact future injuries. Such an inclusion of factors

supports having both upper and lower half data in a holistic analysis. (Garrison et al., 2012)

examined only the shoulder range of motion (excluded hip range of motion like the previous

study) and the impact it has on baseball players with injury to the ulnar collateral ligament.

Simply, this study showed that range of motion deficits in the shoulder were found to be

associated with a tear (grade of tear not explicitly specified) in the ulnar collateral ligament in the

throwing elbow of the examined athlete. This distinction and finding holds important because it

emphasizes a need to include more factors and metrics impacting the arm, as opposed to the

elbow itself. (Garrison et al., 2013) examines the imbalance of the lower extremities as they

pertain to UCL health. This furthers the importance of the NewtForce mound data, which is

readily available within the micro data of this analysis. It provides sound research for the

imbalances of the lower half of the body and the prevalence of UCL injuries in pitchers. (Conte

et al., 2015) examines the prevalence of UCL reconstruction in professional baseball players at

the Major and Minor league levels. While a much more simplistic study than the rest, it is

important to note for background knowledge. 25% of Major League pitchers and 15% of minor

league pitchers have had the surgery. It is important that we understand the prevalence of this

injury to understand the severity of the pandemic of arm injuries. This will be important data to

note for the changes over time, to see how the statistics of prevalence have changed from 2015 to

the currently available data in 2025. (Peters et al., 2018) is a more subjective study as it examines

the success of return to throwing in baseball following UCL reconstruction. More evidence

would be necessary for a greater study and to give more options pertaining to generalizability.

For UCL injuries, returning to throwing is a long, difficult rehabilitation process. The return to

the previous level is less prevalent in professional sports as compared to high school and

collegiate athletes. This shows that age likely plays a much bigger role than originally anticipated

and can be included as a potential variable of interest for analysis. This furthers the importance to

include workload (measured in innings pitched) to make more of a sound, holistic analysis. (Jang

et al., 2019) examines the management of UCL reconstruction processes for baseball players

returning to throwing at the time. The similarity the previous study, (Peters et al., 2018), is noted,

as this study examines the return to throwing rates following reconstructive surgery pertaining to

the UCL. It is important for understanding the odds for success following surgery while also

understanding risks involved with reconstruction. The adds to the knowledge regarding the

prevalence of return to play following a UCL injury process. (Mine et al., 2021) provides an

overview and review of the potential risk factors that are involved with both elbow (UCL) and

shoulder injuries in baseball players. This article does not provide sound conclusions within the

results for findings other than limited shoulder ROM. It is important to provide an overview for

potential risk factors, but not to make any generalizations on the specific results for these risk

factors. Further, it shows the importance to include more statistical factors along with deficiency

metrics when determining injury risk. This proves the need for Baseball Savant metrics, the

macro data, along with the NewtForce and ArmCare App, the micro data. This scholarly

discussion's intended utilization is for the importance of understanding potential risk factors yet

not testing those risk factors based on the specific results of the research. Author and MLB

sportswriter Jeff Passan takes a deep dive into the history of arm injuries in baseball players. He

explores the advances in surgery technology, advances in rehab, history of the injury, and

common misconceptions. He tells individual stories of medical miracles and mysteries,

providing in depth knowledge of the prior history and current state of arm injuries in Major

League Baseball. His inclusion of Doctor Neil ElAttrache provides a parallel to one of the

scholarly sources listed above (Passan 2016). The mention of major surgeons is shown through

one of the variables of the Tommy John data set utilized in analysis.

There is significant research regarding individual deficiencies that have been known to be

causes of significant arm injuries in baseball players. From the previous scholarly discussions

examined, there were findings that, with age, remain wholly relevant, somewhat relevant, and

not relevant at all. The scholarly discussion articles are imperative for providing background

knowledge within the domain and establishing a basis for this holistic analytical approach. Most

discussions included research of variables that should be included: NewtForce, ArmCare App,

Baseball Savant. The extent to which these outstanding variables and information should be

included have been established. This research will encapsulate a holistic approach to individual

deficiencies and combine them into one full approach. While individual research of deficiencies

and exploration is important, it is the holistic inclusion of numerous data points and platforms

that will be incorporated into this research, emphasizing depth, complexity, and offering future considerations based on the findings.

**Ethical Considerations:**

While the data included was collected previously and uploaded to various databases, there are still necessary ethical considerations to note. Alpha Baseball has granted exclusive access of NewtForce Mound and ArmCare App data for consideration in this analysis. The nine players included in the initial data collection have given consent for their data and namesake to be utilized within this analysis. However, their names will be kept out of visuals and other data purposes other than for sorting within the specific data set. The NewtForce data of these Alpha Baseball athletes was uploaded to the cloud system at the Alpha's facility in Mason, Ohio. It was shared with the express consent of player development directors and coaches, Mitchell Bault and Gregg Williams. ArmCare App data of these athletes was obtained in screenshots from the directors, which involved uploading the screenshot data into an Excel file. Baseball Savant data is publicly available for personal usage, exploration, and analysis. Even though there are human subjects involved, there is no blinding involved in this analytical process, as the project does not constitute an experiment. There is no specific treatment group, even though there are two groups of players, in theory (MLB players and Alpha athletes). The last ethical consideration pertains to the usage of a publicly available Tommy John list. Credit belongs to Jon Roegele (MLBPlayerAnalys on X) for creating, updating, and outsourcing the sheet. Created on Google Sheets, the Tommy John Surgery List includes specific information relating to every surgery that has been documented by all professional baseball players, dating back to the very first surgery performed on Tommy John. A specific note is made that some players may be missing, as some organizations tend to be more private with their release of specific medical information

pertaining to Tommy John surgeries and other procedures. Further explanation of variables and

data sets will be explored in the Data section.

<div align="center">**Data & Methods**</div>

**Data:**

MACRO DATA:

Within this analysis, there were several platforms from which data was collected and

utilized. Because the analysis and visuals were performed through RStudio, Sean Lahman's

Baseball Database played a key role. While the database contains various datasets and sub

datasets, only the 'People' dataset was utilized. From the People data set, only the full name of

the player, weight, height, batting side, and throwing side were used. Some cleaning was needed,

as it was concluded that only position players and pitchers that played from 2015 on would need

to be included, because the only players for analysis are those in the Statcast Era. The variable

full name had to be created to make merging of other datasets possible.

The Tommy John Surgery list was next to be implemented into the analysis. The

variables of usage were the player's name, the surgery date; team/level/position at the time of

surgery, throwing side, return date to the same level of competition, and the recovery time in

months. This data set also requires significant cleaning for analysis. Because some players and

pitchers have had more than one Tommy John surgery, they belong to more than one row. The

data had to be cleaned to combine rows for players, so that one player had one row instead of up

to three (as three is the maximum number of surgeries for a player in this list).

From here, the Lahman data set was merged with the Tommy John Surgery list, so that all

players had their specific surgery data in a new data set, should it be applicable that they have

had a surgery.

Next, implementation of the Baseball Savant dataset was constructed. This dataset involved the most complex cleaning process, as it was the largest original dataset to be included for analysis. The Baseball Savant website allows for direct downloading of a CSV file for personal analyses. Overall, 37 variables were selected for analysis. The playerID was automatically included as part of the downloading process. This variable was removed, because the player's name is an important model considering individual injury and surgery data. The variables include: player name, year, innings pitcher, throwing hand, number of pitchers, arm angle (available after 2021, related to shoulder health), and then individual pitch metrics and statistics. Those metrics include the average velocity (in MPH) and usage (%'s) of: four-seam fastballs, sliders, changeups, curveballs, sinkers, cutters, splitters, knuckleballs, sweepers, slurves, forkballs, and screwballs. These individual pitches are then broken down into three categories: fastballs, breaking balls, and off-speed. Further, these average velocities and usage (MPH and %) were calculated for these categories. These categories have been established because different pitchers throw different pitches, some throw multiple, some throw only one or two. Including every pitch that is thrown by pitchers does allow for analysis, and due to the large N of players (8,264) allows for significance for each individual pitch, as well as the three major categories.

MICRO DATA:

The data obtained from Alpha Baseball includes some inconsistent space and time between the different inputs for some individuals. This is not unexpected, as throwing schedules may be misaligned with tournament schedules, lack of logged sessions or visits to the facility, or injuries. The original factors considered pertained only to the NewtForce data collected in each throwing session, however these metrics alone were not enough to provide a significant analysis

of the micro data. Mound data deals primarily with the force output from the lower body and

how it corresponds with different periods of the pitching motion (Newtforce 2025). ArmCare

App measures the force and strength from the upper half, specifically the throwing arm

(ArmCare 2025). NewtForce data includes the force put into the ground towards home plate,

straight down, from back leg to front leg, stride ratio, and player velocity. For the ArmCare App

data, following deletion and mishandling of data, there is only the average strength weight

included along with the player's name. Through this, I have seen that there are other factors to

consider. Height, weight, and handedness (right/left) are relevant inclusions. In addition, BMI is

a variable that can be calculated through height and weight and will be an important predictor for

significance of velocity and injury prediction, with potential for parallel trends with the macro

data. While some of the NewtForce variables will not be utilized, due to their overlap for other

variables, they are all necessary to include for potential analysis. The ArmCare app data includes

metrics related to strength as they pertain to recovery methods and how well the arm can produce

force independent of the lower half force generation.

   Each of these data sources is important when considered independently. When utilized

with one another for a holistic analysis, it becomes increasingly important to analyze the correct

variables and understand how they complement one another. For this reason, the merging of the

datasets was done methodically and carefully to ensure that there were no inconsistencies or

errors when merging. Due to the size of the datasets and the repeated names, it was crucial to

establish effective naming techniques and coding comments, so the correct datasets were utilized

within the individual analyses and data visual creations.

**Methods:**

There is no one best method for utilization, rather multiple statistical tests and methods that will work together, further reference to the holistic analysis, that will aim to complement and validate each other. These include, but are not limited to, multiple regression capturing both logistic and linear methods, time-series analysis, and interaction models to validate statistical significance and provide further evidence of potential relationships. Each of these methods, independently and holistically, will be imperative to assess the mechanical and kinesiological variables and indicators helpful in preventing arm injuries. The regression models intend to show which specific metrics and variables are going to be likely to increase the odds of injury, which are not, and whether these metrics and statistics will be effective in predicting whether a player will become injured or not. The time-series will aim to uncover trends and patterns with the timing of the injury, with the goal of predicting when the injuries may happen again based on metric trends (Bullock et al., 2022). As mentioned within the data and variables sections, the key inputs will include the fatigue/strength/recovery variables, throwing workload, pitch arsenal, height and weight, and then mechanical deficiencies and outliers, finishing up with overall mechanical strengths. The outputs, then, from the logistic regression will be whether there is an injury along with which metrics are sufficient predictors of velocity and injury risk. This also will be similar for the time-series and interaction models (Karnuta et al., 2020). Specifically, the time series will show trends over time, which will give a better overall understanding of when an injury may occur, not simply whether it will in fact occur. A significant portion of variable selection will be based on the method of correlation tests. It will be important to ensure that the variables are not significantly correlated with one another, as that could lead to biased results of the logistic regression. There is significant data for each individual when it comes to the
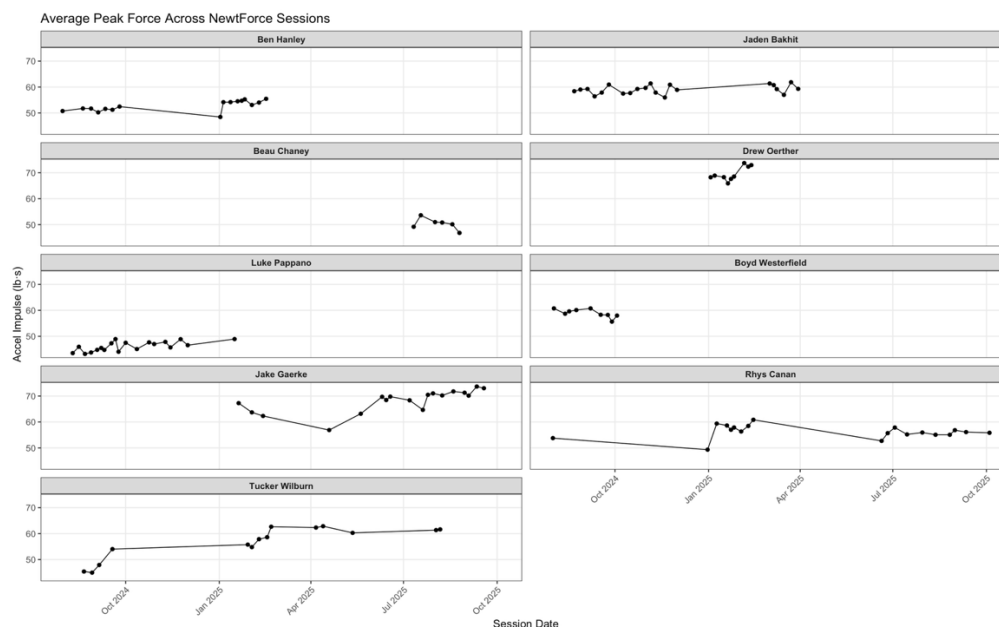
NewtForce dataset (micro), but when attempting to generalize it may be especially useful to add a penalty to account for the smaller sample size. While this analysis will be beneficial as it attempts to uncover specific associations, there will also be potential patterns uncovered that could be useful in prevention of future injuries. Ensuring these assumption and validation steps is important because the conclusions then become more generalizable and practical for their real-world usage, even if there is some inherent issue with the lack of a large N of NewtForce individuals. The size of N will not be a factor of worry for the Baseball Savant and surgery datasets. Overall, this design provides a fundamental basis for the initial analysis and findings of the project.

The order of the methods will follow a specific order. First, it is important to explore the relationship between peak velocity and time. This will be done through the time series analysis that shows the change over time and shows a visual change over time for each individual pitcher in the NewtForce dataset. The time series will give an idea of how the metrics are changing over time through the individual and will be generalizable across the collection of individuals (Bullock et al., 2022). From here, it is sufficient to explore relationships between the important variables within the expanded Baseball Savant dataset. This will show the correlations, through the usage of a heatmap to exemplify the most important variables and their relationships with one another. From here, an integrated comparison will be utilized. Because there is not an exclusive ability to compare NewtForce metrics with Baseball Savant metrics, the few overlapping variables will be compared to show their relationships. This will be completed with scatterplots and best fit lines. The final method explored will utilize a LASSO logistic regression model. This will take the most qualified variables and put them into a model to see how likely a prediction can be made on whether a pitcher will need surgery based on the selected variables.
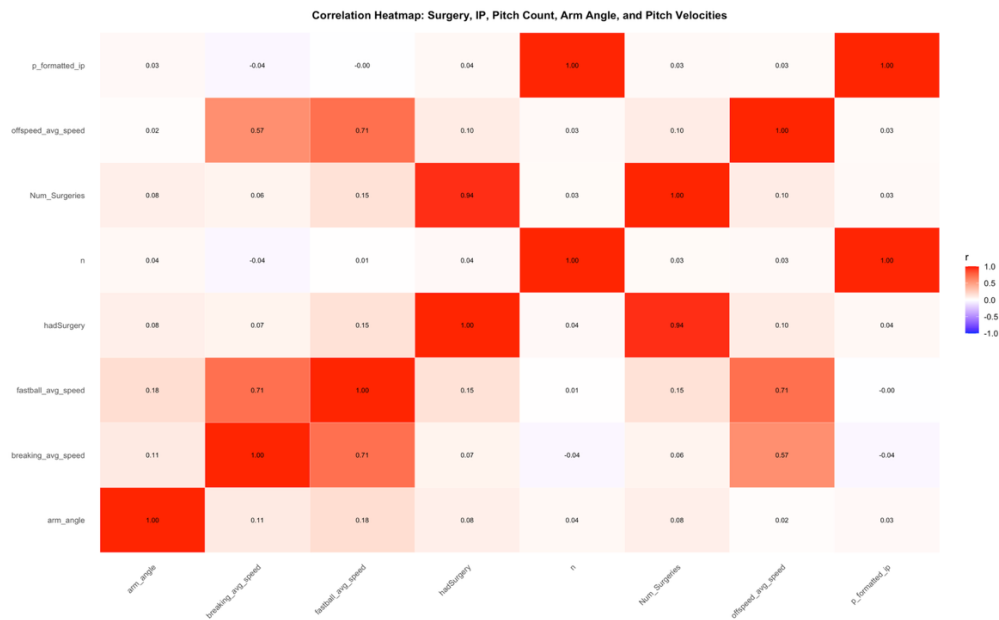
The LASSO model will likely be able to perform better than regular logistic regression due to the presence of more inconsistent trends/results from athletes who do not have consistent data inputs (Karnuta et al., 2020). This methodological approach will provide a sound conclusion to the research question, even if the answer is a null. This holistic approach with methods included will show trends and relationships between various datasets, metrics, statistics, and results.
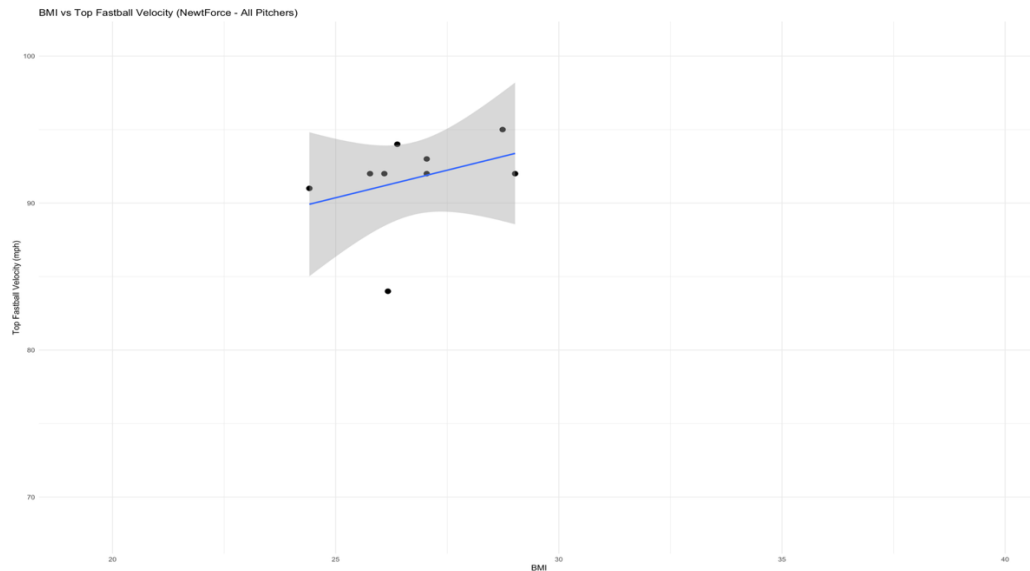
## Results and Interpretation

The results section explores the question: can arm injuries, specifically injuries to the UCL, be prevented or predicted through analysis of various data points and systems that expose indicators of weakness or potential injuries? It is hypothesized that there will be significant relationships between the metrics related to upper half and lower half data. It is expected that each of the platforms of data—Newtforce Mound, ArmCare App, Baseball Savant, Tommy John Surgery, and Lahman—will complement one another when datasets are merged and analyzed collectively. While each platform provides meaningful results independently, it is expected, through integration of the datasets, that trends will emerge that give insight into deficiencies and statistical overlaps that may lead to future injury prediction and prevention. The data sources utilized in this research are metrics used from NewtForce Mound, ArmCare App, Baseball Savant, an independently sourced dataset of Tommy John surgeries, and Sean Lahman's Baseball database. The results follow the order discussed from the method section.

**Figure 1**



The time series in Figure 1 shows the difference in the average peak force across the different NewtForce bullpen sessions for each individual athlete at Alpha Baseball. The average peak force measures the lower half efficiency across players coming back from Tommy John surgery. These measures show growth, yet they also show areas for potential weakness. The peaks show signs of increased strength and positive lower half efficiency, indicating less stress on the arm. The valleys, however, raise questions of whether too much stress is being placed on the arm. These time series do exemplify the changes shown from session to session and attempt to give greater insight to the day and instance of injury during the return to throwing process. It is now sufficient to examine variable relationships of the Baseball Savant data.

**Figure 2**



Correlation Heatmap: Surgery, IP, Pitch Count, Arm Angle, and Pitch Velocities

The correlation heat map in Figure 2 exemplifies the relationships between the key variables in the Baseball Savant dataset. For readability and consistency, only the major variables were included. The variables included fastball, off-speed, and breaking ball average velocities, along with total number of pitches, arm angle, and whether a pitcher had surgery or not, and if so, how many surgeries they had. Very few significant correlations were found, except between velocity variables. These trends would be expected, as the speed of fastball, off-speed, and breaking ball pitches would all trend downward in the respective order for each pitcher, highlighting their strong positive relationship. On the other hand, what would typically be considered a weak relationship, with a correlation of around 0.10, is more important in this case. This includes the relationship and correlation of number of surgeries and average fastball and off-speed velocities. These correlations are important for consideration in the regression model shown later in the analysis.

**Figure 3A**



The bivariate scatterplots shown in Figures 3A, 3B, and 3C all show the relationship between the BMI, body mass index, and fastball velocities. For figure 3A, the relationship does not have a large enough value of N to be considered with any regard to statistical significance. However, it does show a positive trend between BMI and peak fastball velocity. This group of nine pitchers was utilized from the NewtForce dataset obtained from Alpha Baseball. Each of these individuals has had at least once instance of Tommy John surgery in their career. This strengthens support for utilization of the Baseball Savant data as a comparison.
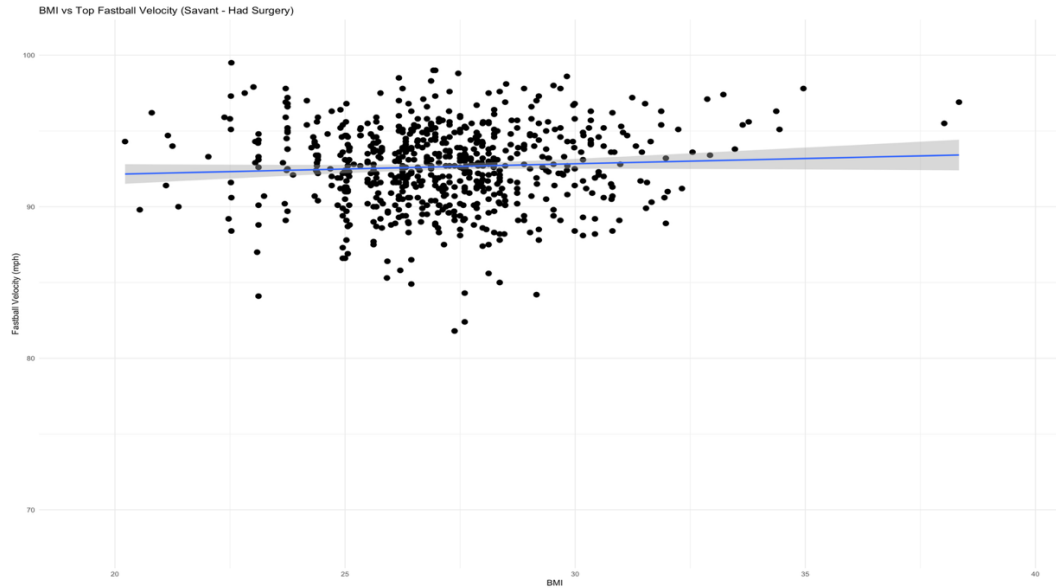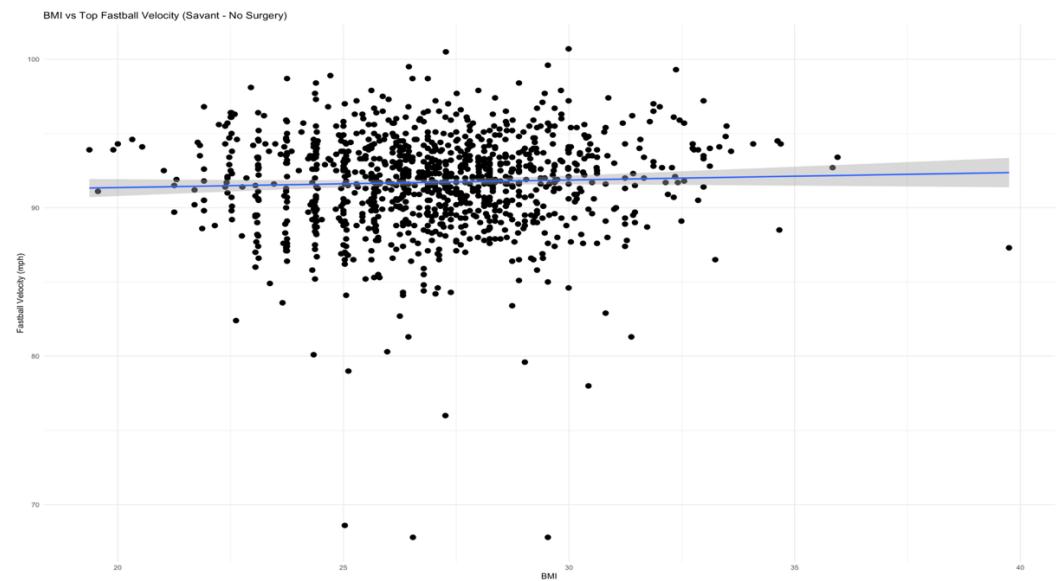
**Figure 3B**



**Figure 3C**



Figure 3B shows the relationship between average fastball velocities of pitchers in MLB who have had at least one instance of Tommy John surgery and their BMI. Figure 3C, on the other hand, shows the relationship between average fastball velocities of pitchers in MLB who have *not* had an instance of Tommy John surgery and their BMI. It is significant to include the analysis of these figures together, because they directly compare the results of two groups: had

surgery and have not had surgery. While there is a greater number of players who have not had

surgery, as compared to those who have surgery, the trend lines appear to show no visual

difference. This further emphasizes the need to examine the individual differences between the

models, which is an interaction model comparing the two scatterplots statistically.

**Figure 4**

```
Call:
lm(formula = velo ~ BMI * group, data = Combined_groups)

Residuals:
    Min      1Q  Median      3Q     Max
-24.044  -1.730   0.253   1.973   8.832

Coefficients:
                            Estimate Std. Error t value Pr(>|t|)
(Intercept)                  71.6099    19.9707   3.586 0.000345 ***
BMI                           0.7501     0.7459   1.006 0.314742
groupSavant_NoSurgery        18.7260    19.9949   0.937 0.349123
groupSavant_Surgery          19.1501    20.0194   0.957 0.338912
BMI:groupSavant_NoSurgery    -0.6990     0.7468  -0.936 0.349386
BMI:groupSavant_Surgery      -0.6809     0.7476  -0.911 0.362568
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.041 on 1761 degrees of freedom
Multiple R-squared:  0.02384,   Adjusted R-squared:  0.02107
F-statistic: 8.603 on 5 and 1761 DF,  p-value: 4.595e-08
```

Figure 4 shows an interaction model comparing the scatterplots shown in Figures 3B and

3C. From the results of the model, we see that BMI is responsible for roughly 2% of all variables

in predicting the average fastball velocity. This shows that BMI is not a significant predictor of

fastball velocity. Further, it shows there is not a statistically significant difference in the average

fastball velocity for MLB pitchers who have had Tommy John surgery as compared to MLB

pitchers who have not had the surgery. The p-value meets the criteria of falling below the 0.05

level of significance, indicating these results are statistically significant.
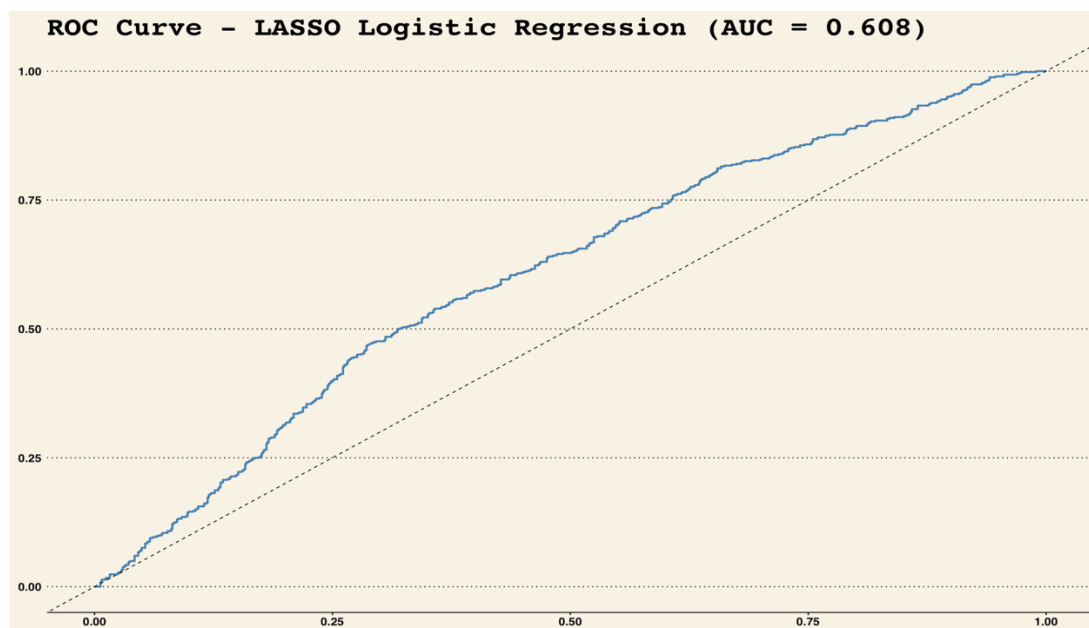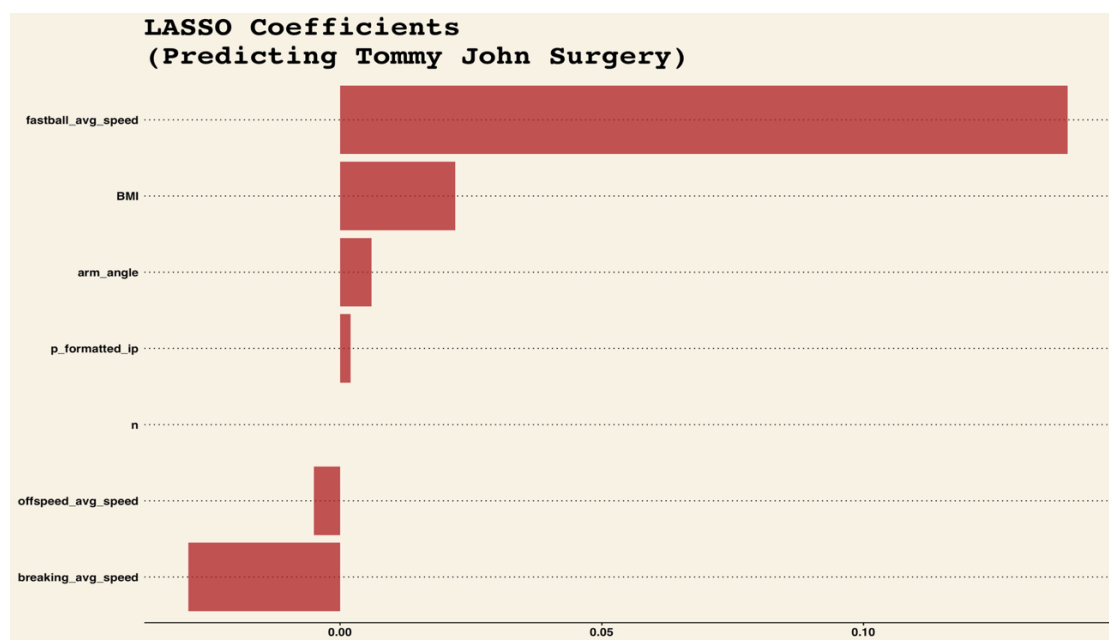
**Figure 5A**



**Figure 5B**

**Figure 5C**

| LASSO Coefficients and Model AUC | | |
|---|---|---|
| variable | coefficient | AUC |
| fastball_avg_speed | 0.139 | 0.608 |
| breaking_avg_speed | −0.029 | 0.608 |
| BMI | 0.022 | 0.608 |
| arm_angle | 0.006 | 0.608 |
| offspeed_avg_speed | −0.005 | 0.608 |
| p_formatted_ip | 0.002 | 0.608 |
| n | 0.000 | 0.608 |

Figures 5A, 5B, and 5C are all related to the LASSO regression model that was created. The goal of the LASSO model was to compile key variables, metrics, and statistics to see if a model could be created that would be sufficient at predicting whether an MLB pitcher would require surgery. Though it would not be able to predict exactly when, it would be sufficient to utilize various key indicators that could show significance when it comes to prevention and prediction of future elbow injury that could lead to a Tommy John surgery. Figure 5A shows how accurate the model is predicting injury prediction. The AUC, or area under the curve, shown in Figure 5C, indicates that the model is average or modest in its prediction of a need for surgery. While this is not an extremely accurate model, it is not a poor model. The variables also indicated in Figure 5C do not indicate extreme impacts. These impacts are exemplified in Figure 5B. It is observed that fastball average velocity is the most significant indicator of a need for surgery, with a coefficient of .139. On the other hand, the average breaking ball and off-speed velocities seem to show adverse results, being indicators of reverting the need for surgery. Other

factors, such as the number of pitches, seem to show no relationship in predicting need for surgery or not. While this model is not extremely efficient, it is not extremely inefficient and shows a modest ability to predict whether an MLB pitcher will, at some point, require Tommy John surgery at some point in his career based on the variables included in the model.

**Discussion**

The important discussion topics are important indicators and inclusions of future analyses, as well as further interpretations of the models and figures. The figures shown were all obtained from an R-Markdown file utilized through RStudio. In the future, and for a final draft, the visuals should be cleaned up and altered for their reproducibility. More ArmCare app data should be obtained for further generalizations to be made regarding the Alpha Baseball athletes. Stride data of MLB pitchers, or lower-half data in some regard would be beneficial for holistic analysis of the MLB data. Greater, more concrete upper half pitching data would be equally beneficial for the Alpha Baseball athletes. If Trackman data could be linked to the NewtForce data, so that individual pitch metrics could be merged for Alpha athletes, the lesser N would be sufficient due to an excess in individual data within each row. The models and figures explored are sufficient, yet not perfect indicators of the relationships and statistics exemplified in the data. More complex, time consuming, and data consuming visuals would be more sufficient for greater generalizability and future analysis. The preceding analysis is holistic, yet there is not insufficient ability to produce further, greater, more complex and time-consuming analysis.

**Conclusion**

The research question for this analysis is, can arm injuries, specifically injuries to the UCL, be prevented or predicted through analysis of various data points and systems that expose indicators of weakness or potential injuries? While the results were not of extreme statistical

significance, the 'null' results do express an answer. Arm injuries are subjective and pertain to

the individual. Because each pitcher has individual differences in genetic makeup, mechanics,

and pitch arsenals, it is difficult to generalize results and findings to groups of pitchers. It does

appear that there are weak associations between faster pitching velocities of fastballs, innings

pitched, arm angle, and BMI leading to a greater chance of requiring surgery from a UCL injury.

It should be mentioned again that these associations are, however, weak. Like other studies,

methods, and results, further analysis of biomechanical data must be observed. It would be

helpful to obtain greater NewtForce results for lower-half metric generalizations to be concrete.

It would also be helpful for MLB pitchers to have lower-half metric data readily available, but

MLB pitcher access to NewtForce mound across the league is impractical. From here, further

methods and results with statistics must be established for further bases and conclusions. These

weak relationships should not be discredited because they are weak, rather understood for their

statistical meaning. Pitchers should proceed with caution and take extra recovery and

rehabilitation methods when throwing consistently at higher velocities in greater volumes.

References

Bartlett, L. R., Storey, M. D., & Simons, B. D. (1989). Measurement of upper extremity torque

production and its relationship to throwing speed in the competitive athlete. The

American Journal of Sports Medicine, 17(1), 89–91.

https://doi.org/10.1177/036354658901700115

Brown, L. P., Niehues, S. L., Harrah, A., Yavorsky, P., & Hirshman, H. P. (1988). Upper

Extremity Range of Motion and Isokinetic Strength of the Internal and External Shoulder

Rotators in Major League Baseball Players. The American Journal of Sports Medicine,

16(6), 577–585. https://doi.org/10.1177/036354658801600604

Conte, S. A., Fleisig, G. S., Dines, J. S., Wilk, K. E., Aune, K. T., Patterson-Flynn, N., &

ElAttrache, N. (2015). Prevalence of Ulnar Collateral Ligament Surgery in Professional

Baseball Players. The American Journal of Sports Medicine, 43(7), 1764–1769.

https://doi.org/10.1177/0363546515580792

Garrison, J. C., Arnold, A., Macko, M. J., & Conway, J. E. (2013). Baseball Players Diagnosed

With Ulnar Collateral Ligament Tears Demonstrate Decreased Balance Compared to

Healthy Controls. Journal of Orthopaedic & Sports Physical Therapy, 43(10), 752–758.

https://doi.org/10.2519/jospt.2013.4680

Garrison, J. C., Cole, M. A., Conway, J. E., Macko, M. J., Thigpen, C., & Shanley, E. (2012).

Shoulder Range of Motion Deficits in Baseball Players With an Ulnar Collateral

Ligament Tear. The American Journal of Sports Medicine, 40(11), 2597–2603.

https://doi.org/10.1177/0363546512459175

Jang, S.-H. (2019). Management of Ulnar Collateral Ligament Injuries in Overhead Athletes.

Clinics in Shoulder and Elbow, 22(4), 235–240.

https://doi.org/10.5397/cise.2019.22.4.235

Major League Baseball. (n.d.). *Statcast | Glossary*. MLB.com.

https://www.mlb.com/glossary/statcast

Mine, K., Milanese, S., Jones, M. A., Saunders, S., & Onofrio, B. (2021). Risk Factors of

Shoulder and Elbow Injuries in Baseball: A Scoping Review of 3 Types of Evidence.

Orthopaedic Journal of Sports Medicine, 9(12), 232596712110646.

https://doi.org/10.1177/23259671211064645

Passan, J. (2016). The Arm: Inside the Billion-Dollar Mystery of the Most Valuable Commodity

in Sports. HarperCollins.

Peters, S. D., Bullock, G. S., Goode, A. P., Garrigues, G. E., Ruch, D. S., & Reiman, M. P.

(2018). The success of return to sport after ulnar collateral ligament injury in baseball: a

systematic review and meta-analysis. Journal of Shoulder and Elbow Surgery, 27(3),

561–571. https://doi.org/10.1016/j.jse.2017.12.003

Petersen, G. (2025, November 6). *UCL injuries rising among younger baseball players, doctors

say*. Https://Www.wmbfnews.com; WMBF.

https://www.wmbfnews.com/2025/11/06/ucl-injuries-rising-among-younger-baseball-

players-doctors-say/

Roegele, J. (Accessed November 2025). *Tommy John Surgery List*. Google Sheets.

https://docs.google.com/spreadsheets/d/1gQujXQQGOVNaiuwSN680Hq-FDVsCwvN-

3AazykOBON0/edit?gid=0#gid=0

Scher, S., Anderson, K., Weber, N., Bajorek, J., Rand, K., & Bey, M. J. (2010). Associations

Among Hip and Shoulder Range of Motion and Shoulder Injury in Professional Baseball

Players. Journal of Athletic Training, 45(2), 191–197.

https://doi.org/10.4085/1062-6050-45.2.191

Ulnar Collateral Ligament (UCL) Injury. (2022). www.nationwidechildrens.org.

https://www.nationwidechildrens.org/conditions/ulnar-collateral-ligament-injury

Wilk, K. E., Meister, K., & Andrews, J. R. (2002). Current Concepts in the Rehabilitation of the

Overhead Throwing Athlete. The American Journal of Sports Medicine, 30(1), 136–151.

https://doi.org/10.1177/03635465020300011201

Appendix

https://github.com/willrettig36/Rettig-DA-First-Draft-401