

Sebastian Cerna



## Business Understanding:

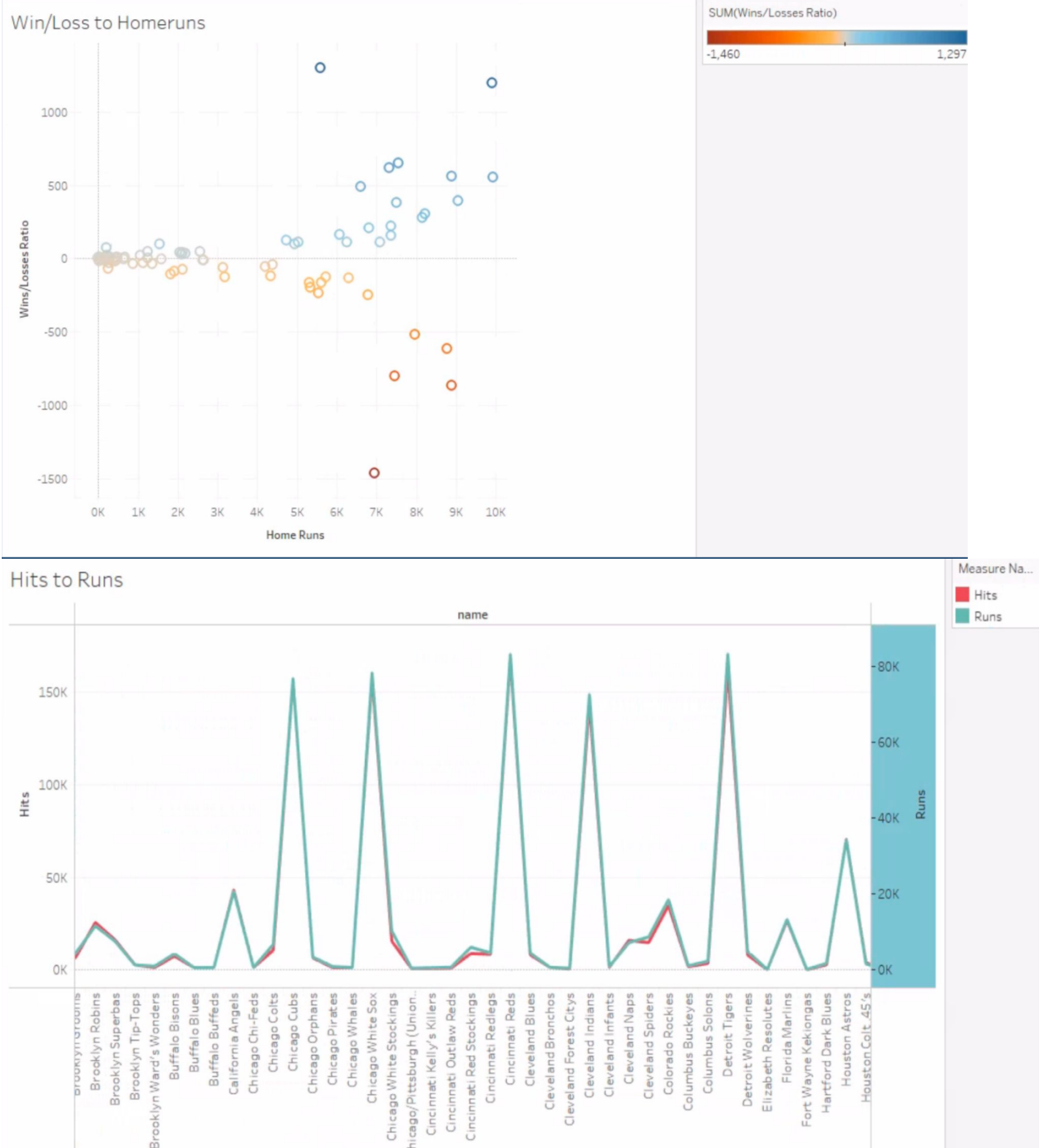
Numerous questions can be answered by our dataset, such as: Which team statistics have the strongest correlation with wins and losses? What impact do statistics like batting averages, home runs, or pitches have on a team's overall performance? Which player qualities are most important for winning games? This dataset's significance is in its comprehensive analysis of team performance, which can help players and coaches make strategic decisions both on and off the field. By focusing on power hitters and strong pitchers, finding metrics which produce the best results, and improving recruitment strategies, it assists major league baseball organizations in improving their team structure. Coaches can better allocate resources, prioritize training programs, and adjust game plans for the most wins by examining the connections and trends within each team.

## Data Understanding:

Using K-means clustering we were able to determine the most important attributes that a team needs to increase their win percentage. In our K-means clustering we created four different clusters using the number of all-star players on a team, the win percentage, runs per season, and strikeouts per season. We noticed that the most important attributes to a team's win percentage was how many Allstars they had. Team with only one all-star had the lowest win percentage, while teams with 5 Allstars had the highest win percentage. The mid-tier teams average 3 Allstars. From the clustering we can also see that the worse performing teams had significantly less runs and strikeouts than the higher performing teams. The mid-tier performing teams were either good at pitching or good at batting while the top teams weren't the best at pitching or hitting, but they had good performance in both. Wins and losses, home runs, pitching metrics, hitting statistics, and all-star players were among the crucial fields we used. A scatter plot illustrating the correlation between home runs and victories is our initial visualization. This graphic shows how home runs may win games on average, but it also illustrates how it varies depending on how many home runs the sides score. A line chart that contrasts a team's hits and runs is the second visualization. The graph makes it clear that there are more runs than hits. The third graph helped us visualize our finding that the amount of Allstars a team has directly correlates to their win percentage. You can see an increase in a team's win percentage every time they gain an all star and their win percentage dips when they lose one. There are some outliers



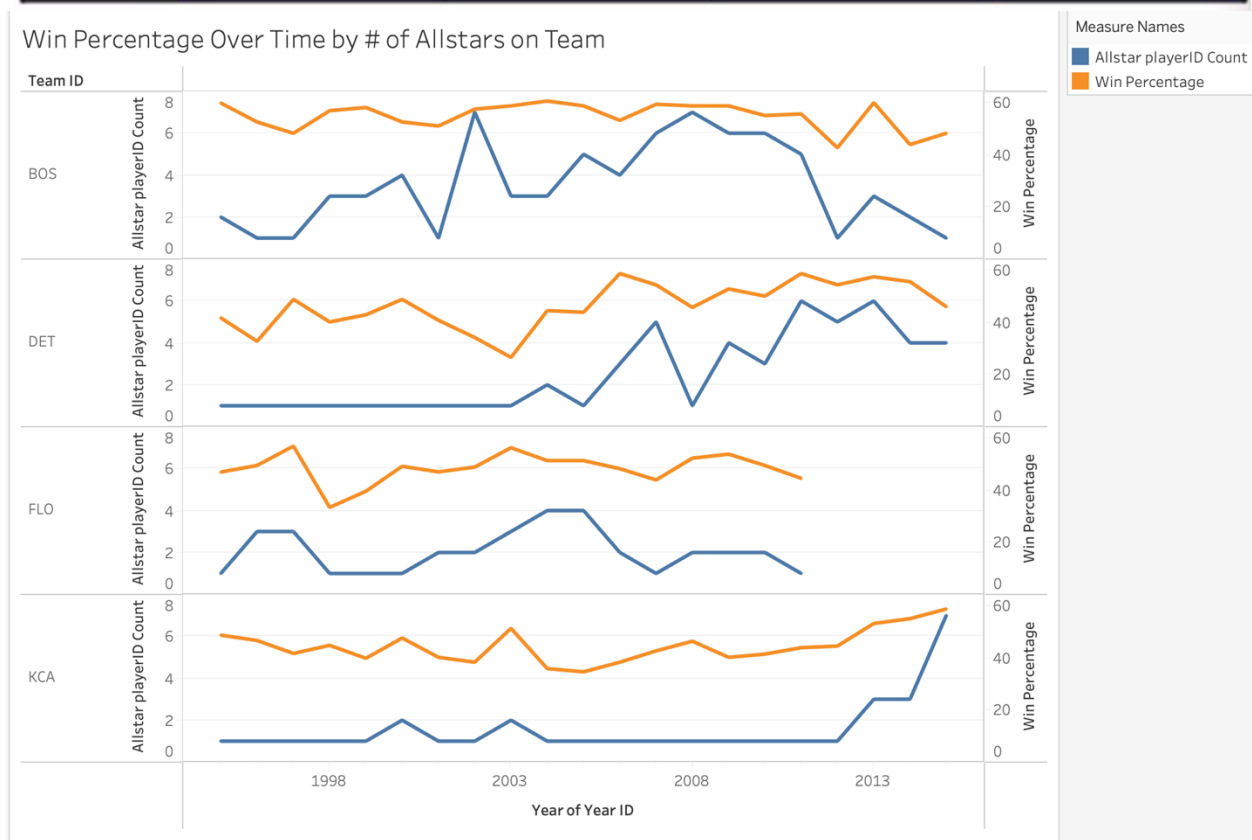
where the win percentage doesn't directly follow the all-star count, but majority of the time it does.





<b>Size</b>	34.5% (215)	10.9% (68)	25.5% (159)	29.1% (181)
<b>Inputs</b>	Allstar_playerID_Count	Allstar_playerID_Count	Allstar_playerID_Count	Allstar_playerID_Count
	WinPercentage 42.93	WinPercentage 58.81	WinPercentage 51.54	WinPercentage 53.68
	R 699.40	R 831.60	R 689.83	R 826.83
	SOA 1,017.08	SOA 1,152.24	SOA 1,229.94	SOA 1,046.91

Win Percentage Over Time by # of Allstars on Team

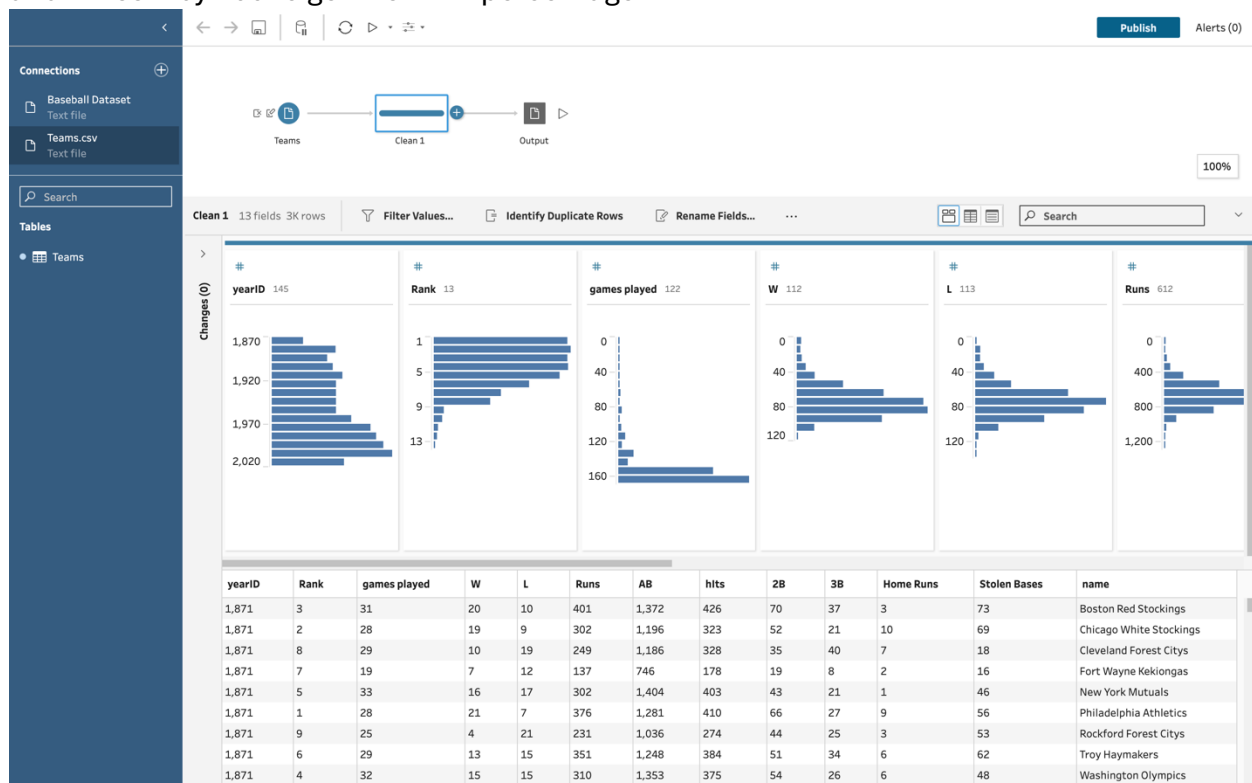


## Data Preparation:

There were more than 60 fields in a single area of the Kaggle dataset, which included 12 distinct datasets in total. We didn't use the management, recruiting, etc.; we just used the



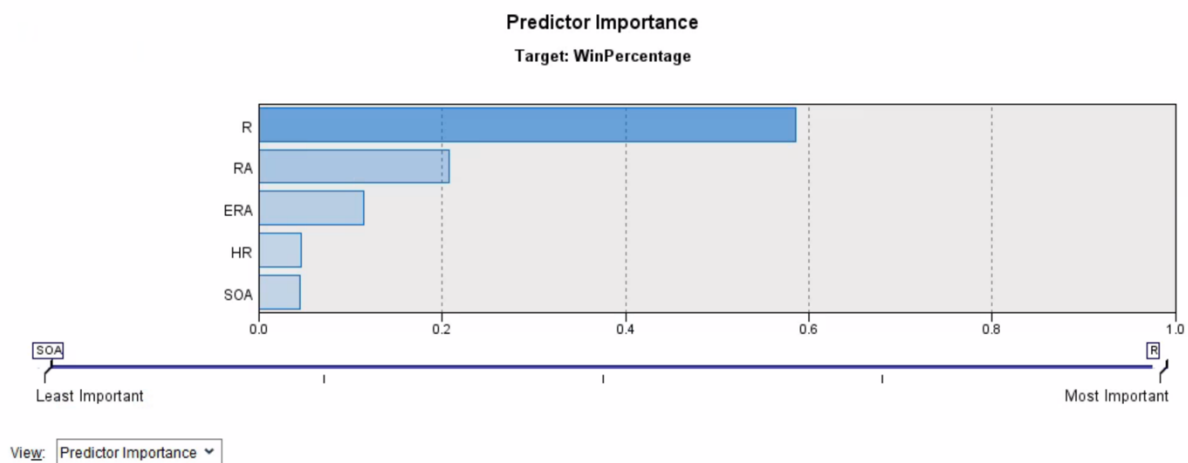
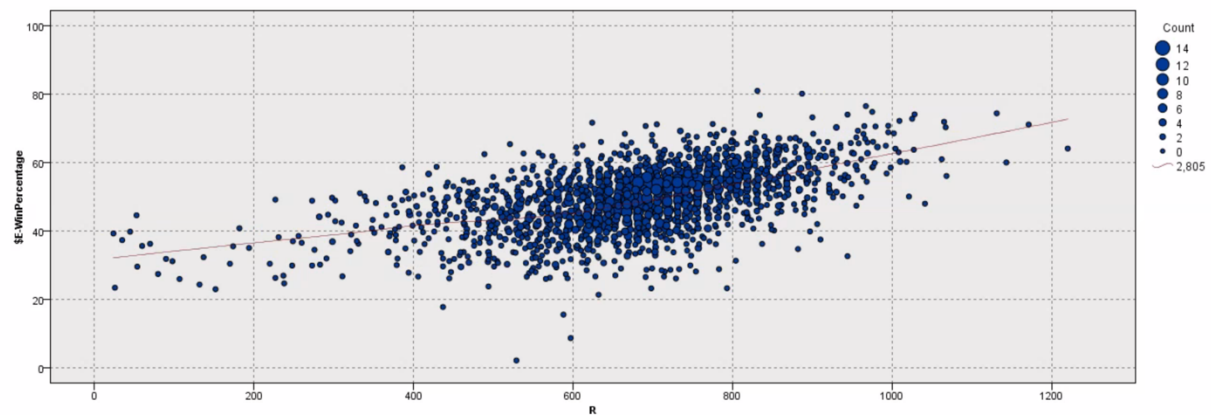
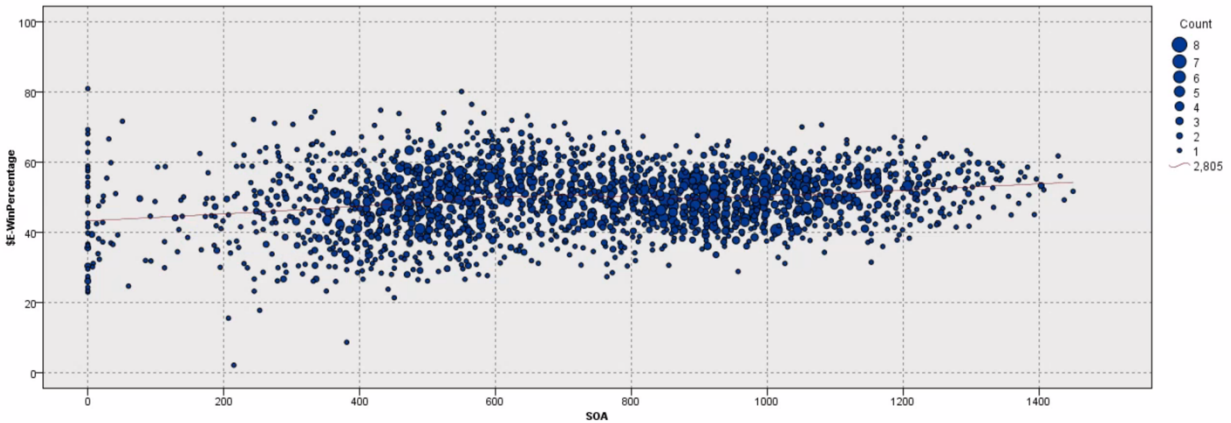
data from each team. Numerous null values, duplicate entries, and data inconsistencies had to be eliminated. We already had the name of the baseball league, so I deleted the team IDs that we didn't need. In our initial visualization, my group identified outliers in the home runs that were won and lost. For the clustering and the predictive models we had to use spss to merge tables from the database together. We also got rid of any data that happened before the allstar game was started. We also created a new field called win percentage that took the amount of wins a team had and divided them by the total games and times it by 100 to get their win percentage.



## Model Building:

In this analysis, a regression model was built to predict a team's win-loss percentage (W-L%) using several key factors such as runs scored (R), home runs (HR), runs allowed (RA), earned run average (ERA), and strikeouts by pitchers (SOA). The goal of the model was to assess the impact of these variables on a team's overall performance and forecast the win percentage. According to the model runs has the greatest impact on win percentage while strikeouts by pitchers (SOA) had the least impact on our expected win percentage. From the graphs that visualize our regression model we can see the win percentage to strikeout slope is not as big as the slope for win percentage and runs. It is also important to notice that HR didn't have as big of an impact on win percentage as we previously thought. For teams it might be more important to focus on getting on base instead of trying to hit homeruns.

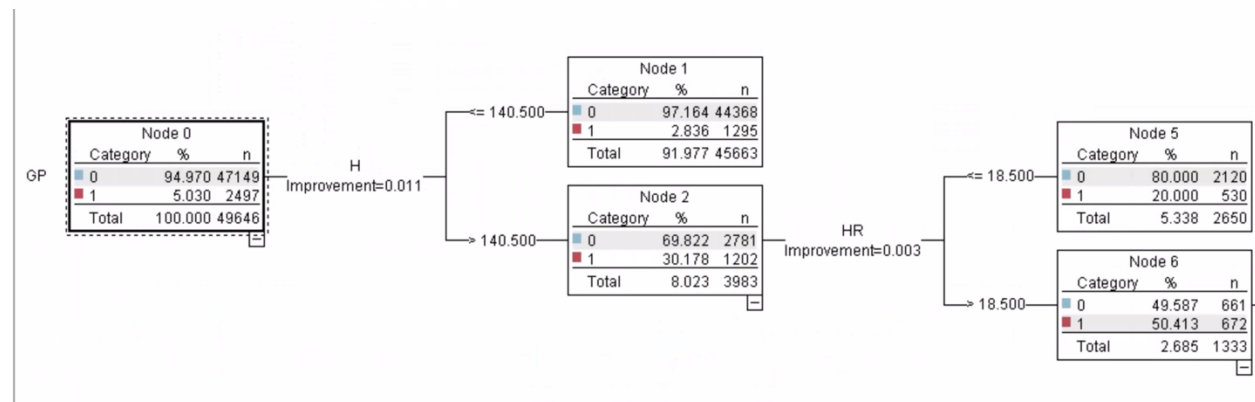
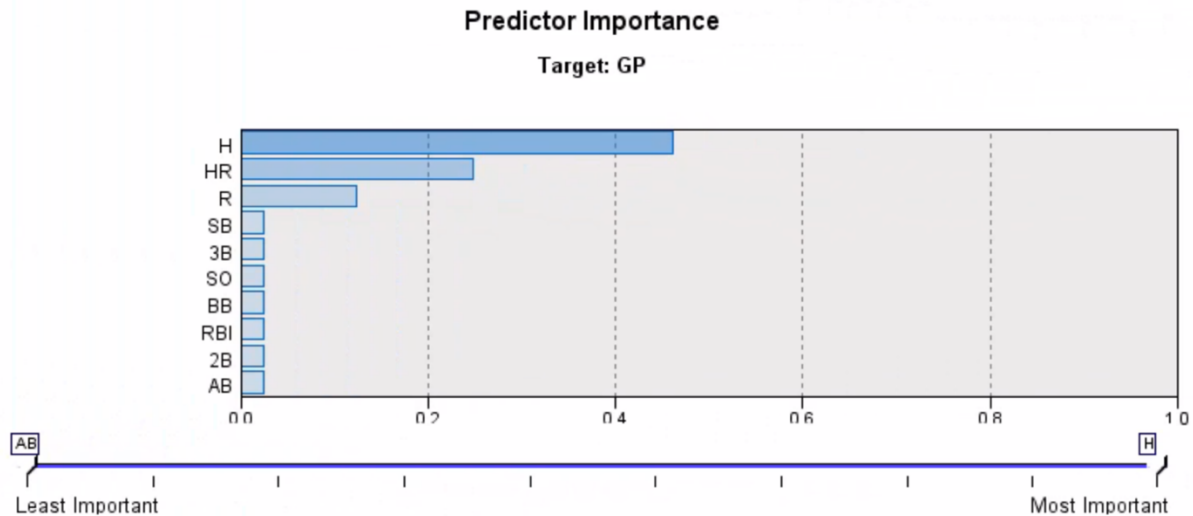




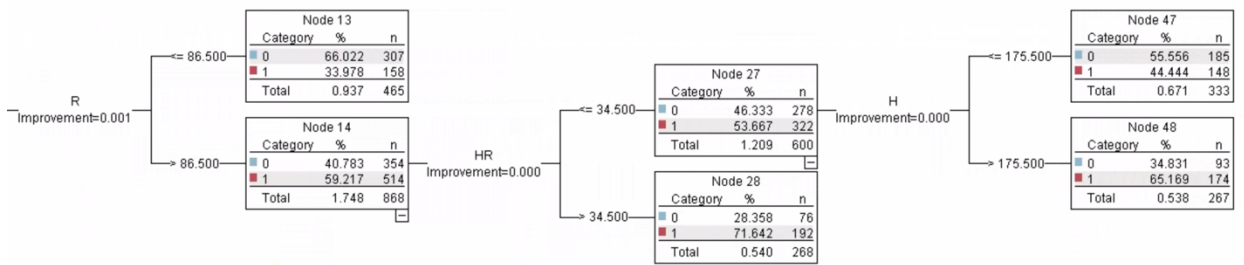
Our second model is a decision tree which we wanted to use to predict which batters would become Allstars based on their season stats. We did this by merging the batters table with the all-star table. We used an outer join which made it so any players that didn't make the all-star game they had a null value in all-star game played and then turned all the null values to 0 using a filler node. We wanted to include all our batting stats into our decision tree and let it decide what the most important stats were. As we can see the most important stats that determine if a player is an all-star is hits, then HR, then runs.



According to our decision tree If a player has fewer than 140.5 hits, the model immediately predicts they will not be an All-Star. Players with more than 140.5 hits and fewer than 18.5 home runs are also predicted to not be All-Stars. Players who hit more than 18.5 home runs and have more than 86.5 runs might be candidates for the All-Star game, depending on their total home runs and hits. We noticed the model was better at predicting if a player wasn't an allstar as opposed to predicting if he was an allstar. We believe this might be because overall in history there have been very little amount of allstars compared to the total player population.







```

H <= 140.500 [ Mode: 0 ] => 0
H > 140.500 [ Mode: 0 ]
  HR <= 18.500 [ Mode: 0 ] => 0
  HR > 18.500 [ Mode: 1 ]
    R <= 86.500 [ Mode: 0 ] => 0
    R > 86.500 [ Mode: 1 ]
      HR <= 34.500 [ Mode: 1 ]
        H <= 175.500 [ Mode: 0 ] => 0
        H > 175.500 [ Mode: 1 ] => 1
      HR > 34.500 [ Mode: 1 ] => 1
  
```