

Do You Really Need Trash Cans?

February 26, 2020

The 2019–2020 MLB off-season was dominated by the revelation that the Houston Astros stole signs by looking at a live camera feed from center field in an area just off the dugout and communicating to predict what pitch was coming next? Do you even need to be stealing signs, or just a computer?"



Photo by Lesly Juarez on Unsplash

The Astros Scheme

The Astros scheme has been very well covered by now. I won't delve into the specifics of their cheating scheme. First, they only seemed to bang the trash can if the next pitch wasn't going to be a fastball. Second is that Rob Arthur from Baseball Prospectus determined the time.

"By and large, the Astros tended to get the signals right, but it was hardly perfect. They were more right than wrong of the time and they were wrong seven percent of the time. ... Based on Adams' data, the Astros were not silent."

(Source: <https://blogs.fangraphs.com/the-most-important-bangs-of-the-astros-scheme/>)

So, in order for our machine learning model to be comparable to the Astros' scheme, it must predict the next pitch.

The Data

The great thing about applying analytics and machine learning to baseball is that there is a wealth of data. I used a Python package, which provides a wrapper for Statcast data, which has entries of every pitch thrown. The count is. Also included are what hand the pitcher throws with and the hitter's stance. Crucially, Statcast has a previous pitch velocity column.

Once the data was pulled from Statcast and massaged a bit, I split it into a training and testing set. Since there were many more cases of pitchers throwing a fastball than there were other types of pitches, I achieved very good accuracy by predicting the next pitch was always going to be a fastball. Obviously, this is not a good model, leaving us with a balanced training dataset between all different types of pitches.

In this analysis, we will use two different datasets, one consisting of Jose Berrios' pitches from the 2019 season and another consisting of all pitches in every MLB game between April 1st, 2019 and April 7th, 2019. These two datasets will allow us to test our model in deciding what pitch to throw in a given situation.

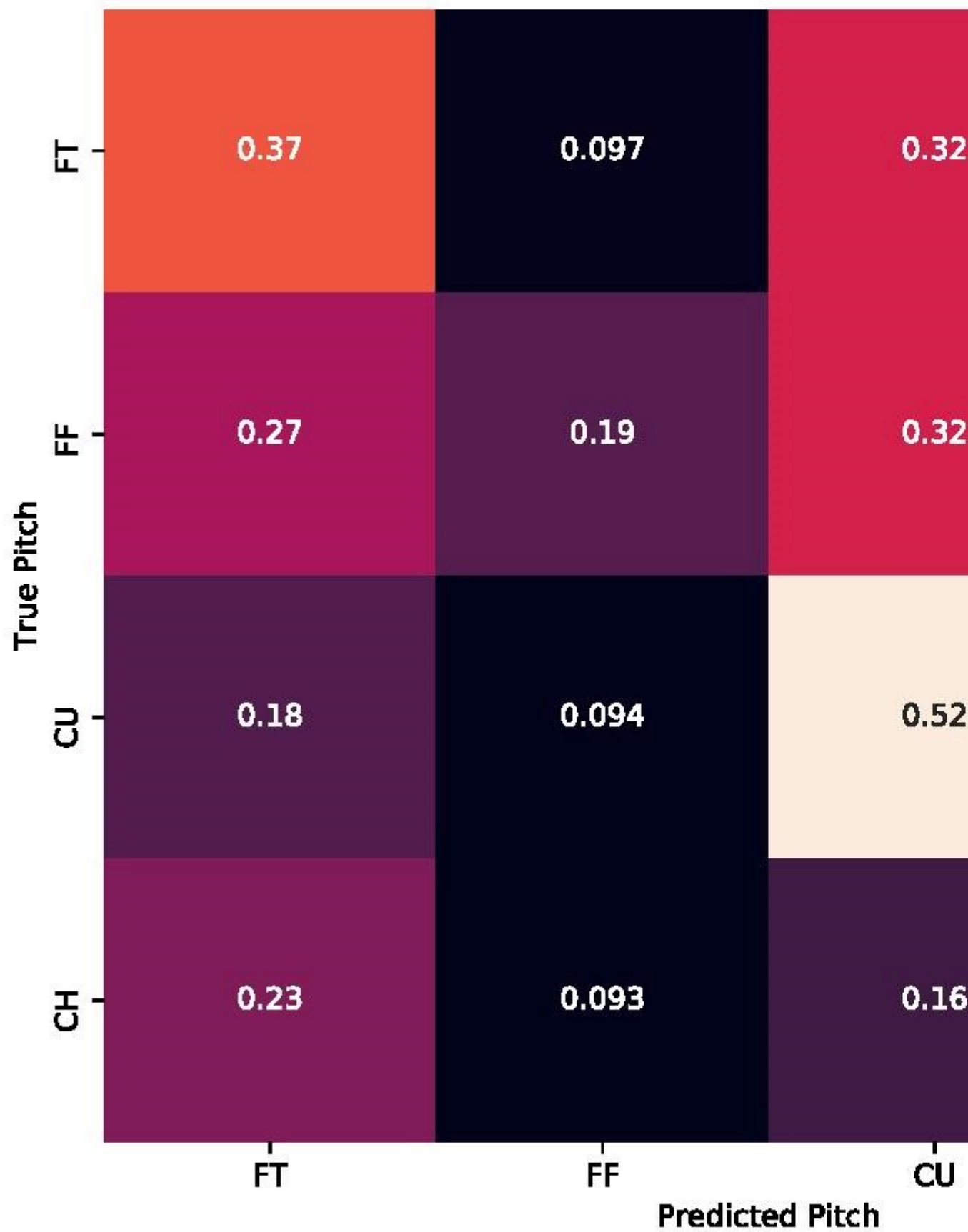
With these two different datasets, I also created a version of each that only contains if the pitch is a fastball or not. I am taking two-seam fastballs, and cutters as "fastballs" and everything else as a non-fastball. I am taking a

The Models

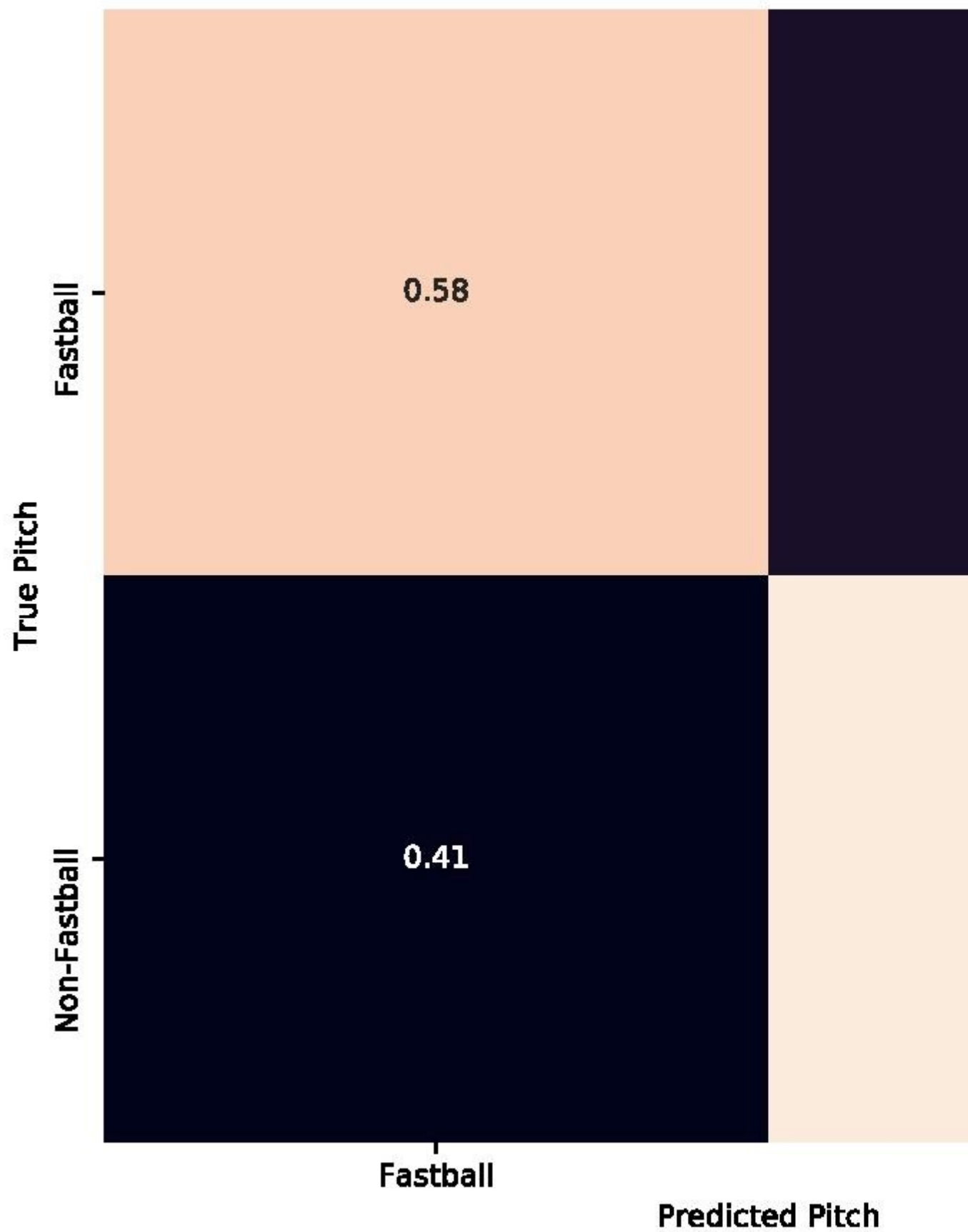
I started with some “basic” machine learning models such as Decision Trees, SVMs and k-Nearest Neighbors. I then decided to feed the data into two different neural networks: one for each version of the dataset.

The Jose Berrios Models

These models achieved better accuracy than our other approaches, with the multi-pitch classifier achieving 95% accuracy. It is unclear to me why there is such a discrepancy between these two classifiers, but I will investigate this further.



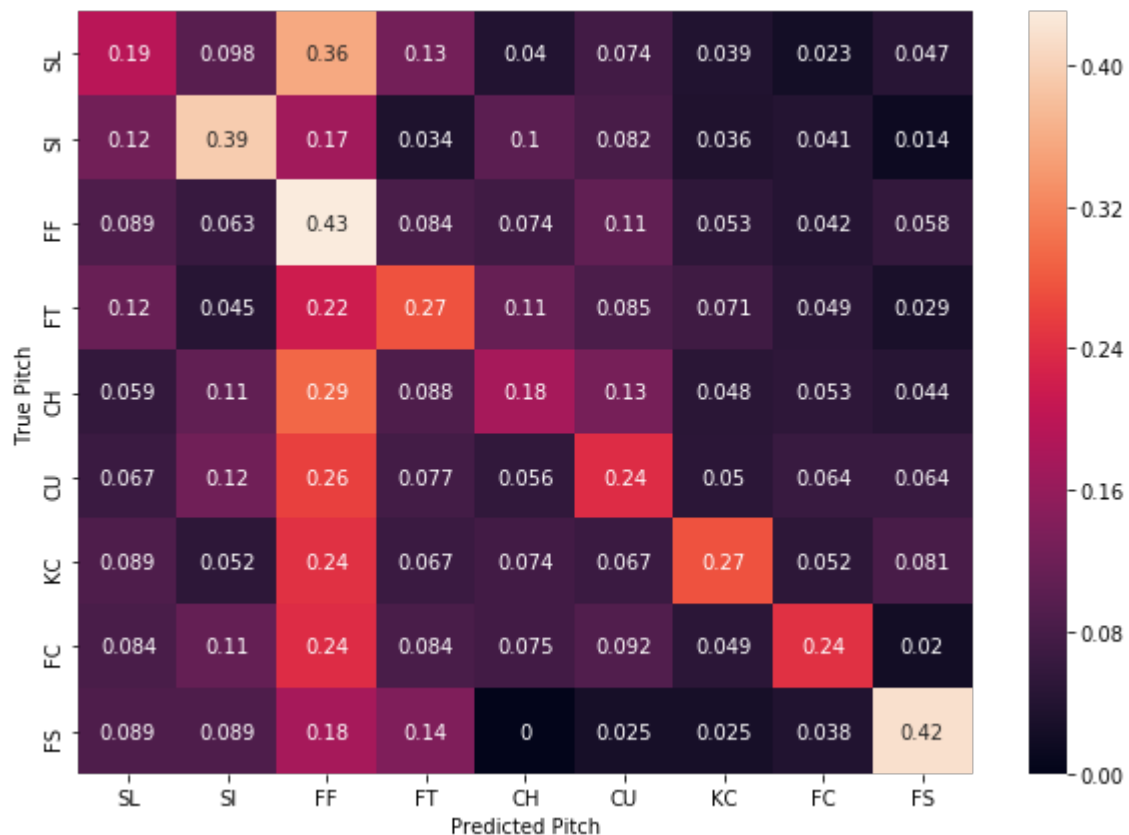
Looking at the confusion matrix above, we can see how many of the pitches of a given type matched the model's prediction. For example, for curveballs (labelled "CU"), we can see that 52% the model predicted successfully, while 20% of the time it predicted at 18% of the time, followed by a four-seam fastball at about 10% of the time. If you turn the accuracy of this model was pretty good, turning its predictions into actionable results would be d



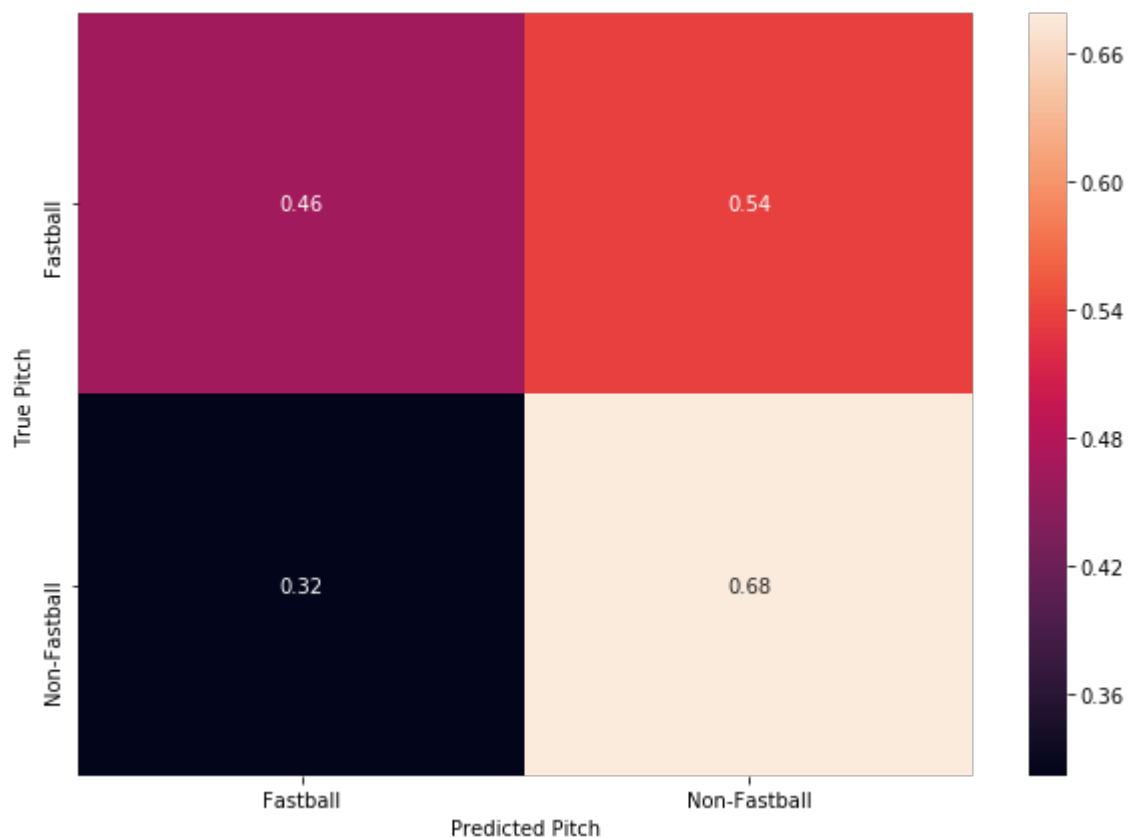
Above is the confusion matrix for the fastball/non-fastball classifier. As you can see, our true positive rate for the multi-pitch classifier, the outcome was relatively balanced. This is a good thing, as we want to recognize the worse outcome for our trash can banger: only about 60% of the time they would be correct. Maybe

The League-Wide Models

The league-wide classifier for all pitches was 88% accurate on the testing set, but as before, let's



As you can see, there isn't that much improvement in terms of improving the trash-can bang's accuracy. The classifier is going to be thrown 43% of the time it actually is, with it predicting that a curveball or a slider is going to be thrown. Correctly predicting that the pitch is going to be either a four or two-seam fastball is only about 50%. Maybe the fastball/non-fastball classifier will produce better results this time?



So, obviously this isn't a great outcome either. The test set accuracy was 57%, and as the confusion matrix shows, the model would bang 54% of the time there was a fastball coming, when it is only supposed to bang on a non-fastball. It also had a false negative rate of 32%. This means that the trash can wouldn't be banged 32% of the time.

Conclusion

Coming into this project, I thought that machine learning might be able to predict what the next pitch would be. Even with some models producing an accuracy of greater than 75%, the false positive and false negative rates were still high. More architectures of neural networks and more data surrounding each pitch could possibly produce better results. But, knowing what the next pitch will be is to cheat.

Checkout the code here: <https://github.com/parkererickson/baseballDataScience/blob/master/n>