

# MLB Pitch Prediction

John F. Adamek

2023-07-31

## Introduction

The use of machine learning models in baseball have become popular thanks to the efficient ‘learning’ capabilities of these models to learn what an outcome (e.g., pitch type) should be based on a number of features (e.g., initial speed, horizontal break) fed into the model. This is considered training the model. The goal is to train a model to be as accurate as possible when predicting an outcome. Here I show a basic example of using general and pitcher-specific machine learning models to predict a pitcher’s pitch type (i.e., two-seam fastball, four-seam fastball, curveball) based on the characteristics of the pitch. The data consists of six different pitchers.

## Machine Learning Layout

Goal: Give the most likely pitch type for all of the pitches in the test dataset using information from the training dataset.

The goal is to predict the type of pitch from the training set by only given a numerical value associated with the pitch type and not the actual name. This will be done through a series of steps:

**Step 1:** Check and visualize the data.

**Step 2:** Prepare the data to be fitted to each of the models.

**Step 3:** Evaluate model performance by examining its accuracy in predicting pitch type in the testing dataset

**Step 4:** Determine the model with the highest accuracy scores to predict pitch type in the testing dataset

**Step 5:** Make final predictions

**Step 6:** Check and visualize the predicted results to the original data. To see if patterns match.

## Methods

**Step 1:** The first step is to look at and visualize the data. What are the variables in the provided dataset? The basic descriptive means of the independent variables and observations for each pitcher were displayed. Findings show that the pitchers in this dataset are likely to be right handed pitchers due to their release point (initposx) being on the third base side of the pitching rubber (Tables 2 and 4). Additionally, we can see that pitch type 9 and 10 are most likely refer to fastballs due to greater initial speed with pitch type 9 associated with a 2-seam fastball/sinker and pitch type 10 associated with a 4-seam fastball based on greater horizontal movement towards a right-handed hitter (breakx) for type 9 and lesser vertical movements downward (breakz) for type 10. Furthermore, Pitcher 3 has only 12 observations (pitches) in the train set which is not an efficient sample size to train and test a model for future predictions. Therefore, I will take this in consideration when determining the model to be used for final predictions. I will test separate models for individual pitchers and the total model performance for addressing Pitcher 3 and Pitcher 6. As expected,

the correlation matrix show's significant ( $p < .05$ ) correlations amongst independent variables ruling out regression based models such as logistic regression.

Based on the data and research question, I will fit and evaluate the performance of three machine learning classification algorithms: decision tree (DT), k-nearest neighbor (K-NN), and support vector machine (SVM).

Table 1: Pitch Classification Dataset

Variables	Description
pitchid	a unique identifier for each pitch
pitcherid	identity of the pitcher (1-6)
yearid	year in which the pitch occurred (1-3)
height (in)	height in inches of the pitcher
initspeed (MPH)	initial speed of the pitch as it leaves the pitcher's hand
breakx (in)	horizontal distance where a pitch crossed the plate in relation to a hypothetical spinless pitch
breakz (in)	vertical distance where a pitch crossed the plate in relation to a hypothetical spinless pitch
initposx (ft)	horizontal position of the release point of the pitch
initposz (ft)	vertical position of the release point of the pitch
extension (ft)	distance in front of the pitching rubber the pitcher releases the ball
spinrate (RPM)	how fast the ball is spinning as it leaves the pitcher's hand
type	type of pitch that was thrown

Table 2: Basic Means of Variables

type	mph	spin	breakx	breakz	initx	initz	ext
2	76.985	2512.840	4.892	-6.841	-1.772	5.952	6.193
3	82.459	1867.164	1.059	-0.033	-1.383	5.754	6.207
4	83.821	1364.294	-5.607	3.453	-1.744	5.895	6.196
7	84.628	988.921	-2.907	2.479	-1.023	5.993	6.206
8	88.903	2346.131	1.233	4.711	-1.815	5.813	6.205
9	91.151	2065.167	-7.126	6.907	-1.828	5.856	6.204
10	92.141	2131.526	-3.185	9.447	-1.627	5.942	6.195

Table 3: Total Observations(pitches) for each Pitcher in Training Set

pitcherid	N
1	1049
2	2137
3	12
4	1840
5	5609

*Note* Pitcher 3 has n=12 observations

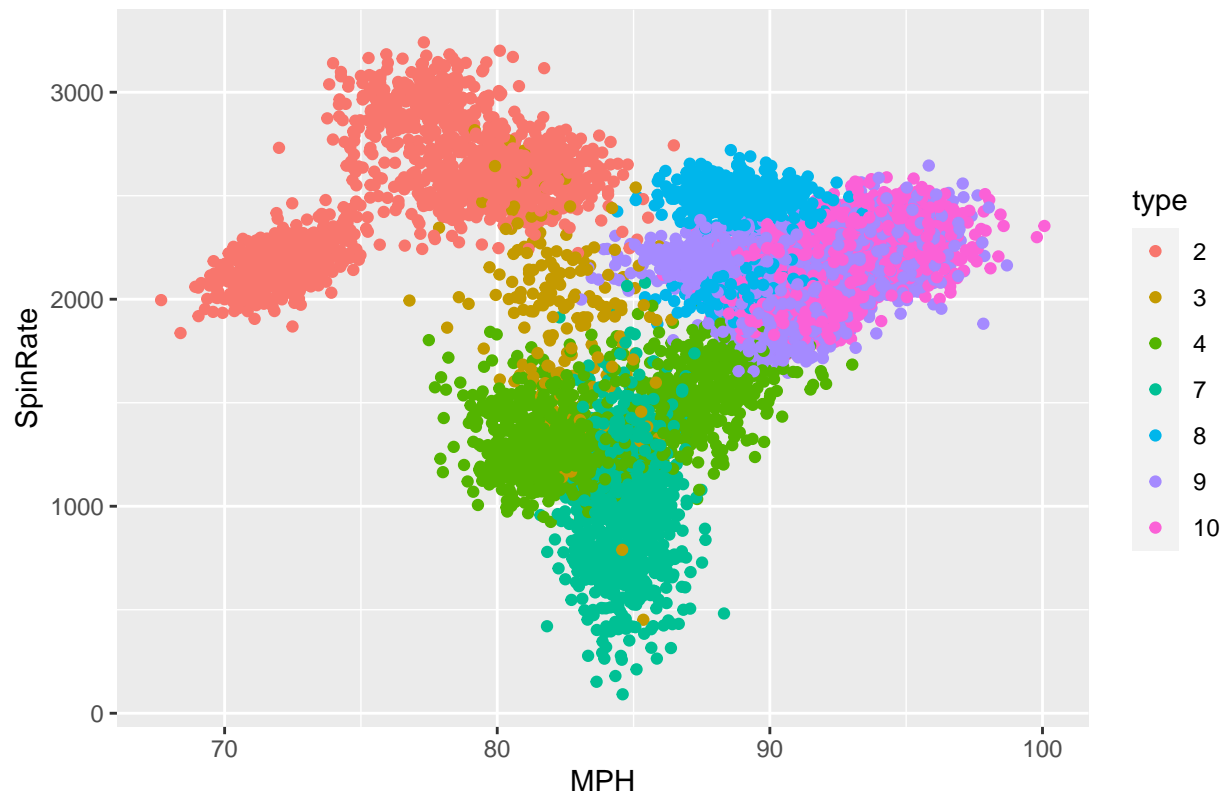
Table 4: Means of Variables for Individual Pitcher by Type

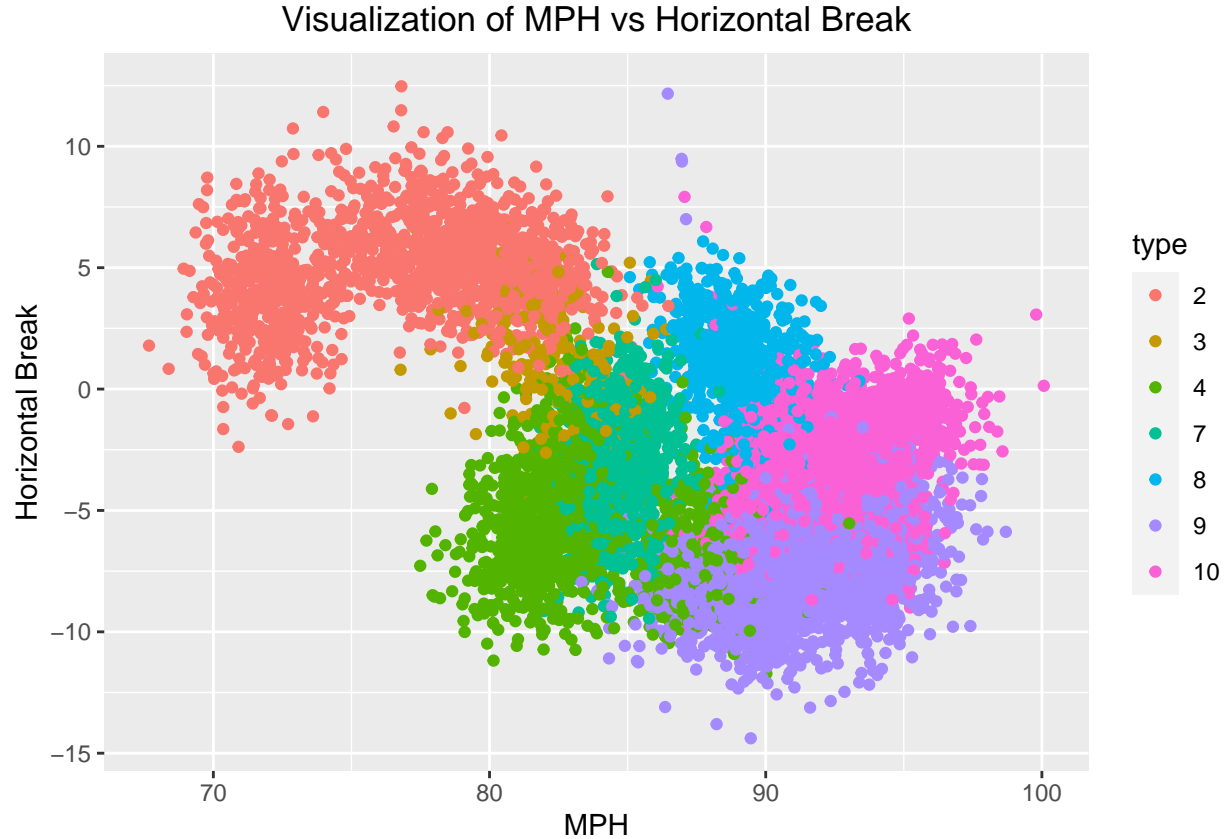
Pitcher	type	mph	spin	breakx	breakz	initx	initz	ext	pitches
Pitcher1	2	77.185	2951.308	5.829	-6.466	-1.854	6.397	6.183	257
	4	81.256	1432.336	-6.688	0.901	-1.838	6.413	6.198	255
	9	88.111	2206.983	-8.489	3.433	-2.062	6.275	6.186	421
	10	89.380	2232.940	-5.992	7.134	-1.886	6.405	6.194	116
Pitcher2	2	79.726	2574.145	5.550	-6.765	-2.184	5.860	6.199	505
	4	87.640	1619.410	-5.492	3.156	-2.352	5.743	6.185	257
	9	93.987	2208.271	-6.268	7.541	-2.265	5.758	6.214	614
	10	93.990	2241.333	-1.666	8.986	-2.268	5.856	6.196	761
Pitcher3	3	84.724	2044.592	-0.095	4.388	3.885	6.662	6.212	2
	9	86.873	2041.951	9.506	4.935	4.145	6.437	6.218	4
	10	87.670	2098.654	4.792	8.584	4.069	6.510	6.136	6
Pitcher4	2	81.272	2624.056	4.391	-2.727	-1.920	5.849	6.196	231
	4	87.025	1430.257	-7.692	3.305	-2.151	5.694	6.183	192
	8	88.950	2490.009	1.946	4.423	-1.990	5.846	6.210	444
	9	93.422	2309.789	-7.346	7.695	-2.076	5.860	6.192	303
	10	93.364	2336.793	-4.694	9.769	-1.968	5.923	6.196	670
Pitcher5	2	72.034	2167.257	3.957	-9.054	-1.235	5.861	6.190	490
	3	82.437	1865.416	1.070	-0.077	-1.435	5.746	6.207	203
	4	82.316	1211.610	-4.574	4.660	-1.331	5.807	6.203	626
	7	84.628	988.921	-2.907	2.479	-1.023	5.993	6.206	901
	8	88.813	2068.383	-0.143	5.268	-1.476	5.751	6.195	230
	9	90.447	1927.528	-7.096	7.429	-1.568	5.782	6.208	1610
	10	90.928	1981.326	-3.100	9.710	-1.167	5.956	6.194	1549

Table 5: Correlation Matrix of Independent Variables

	initspeed	breakx	breakz	initposx	initposz	extension	spinrate
initspeed							
breakx	-0.54***						
breakz	0.87***	-0.60***					
initposx	-0.21***	0.07***	0.02*				
initposz	-0.13***	0.08***	-0.07***	0.21***			
extension	0.01	-0.01	0.01	0.01	-0.01		
spinrate	0.11***	0.37***	-0.10***	-0.45***	0.05***	-0.01	

Visualization of MPH vs Spinrate for Pitch Types





## Data Preparation

**Step 2:** The independent variables were first normalized to ensure the units were properly scaled. Prior to determining which algorithm to use for predicting the final pitch type, the dataset was split (75%/25%) into a training and testing set in order to evaluate model performance for the three different machine learning algorithms. The training set will be used to train each of the models which would then predict pitch type on the testing set. Model performance is evaluated based on the models ability to accurately predict the pitch type in the testing set. In addition, the training set was further separated for each of the five pitchers to run six separate models (five for each pitcher and one with data from all five pitchers) for the DT and K-NN. Models will be evaluated and compared based on their ability to accurately predict pitch type in the testing set. Because Pitcher 6 does not have any data to train on, total model performance will be used to predict pitch type for Pitcher 6 in the testing dataset. Additionally, due to the limited amount of data available for Pitcher 3, I expect to use the total model performance to predict pitch type for Pitcher 3 in the testing dataset as well. If accuracy for the total model is greater then accuracy for the separate models, the total model will be used to predict performance for all pitchers. Otherwise, the individual pitcher data will be used to predict that pitchers pitch type in the testing dataset. For instance, if the K-NN model had a greater predicted pitch type accuracy for Pitcher 2 compared to the total K-NN model then the model for Pitcher 2 will be used to predict pitch type for Pitcher 2 in the testing dataset.

## Models

**Step 3:** For each of the three algorithms, separate models were trained on the training set and then predictions were made on the testing set (with the dependent variable, pitch type, removed). The results from the models predictions were compared to the actual results with performance being represented by an accuracy percentage.

## Decision Tree

Six separate decision tree's were created, five for each pitcher and a total model using data from all five pitchers. After training the data for each model and making predictions on the testing set, the total model performance was 84% accurate in predicting pitch type. Greater performance was found for the separate models for Pitcher 1 (93%), Pitcher 2 (94%), Pitcher 4 (87%), and Pitcher 5 (88%) with an expected low accuracy of 66.67% for Pitcher 3 (Table 5).



Table 6: Decision Tree Model Performance

Model	Accuracy
Total Model	0.84
Pitcher1	0.93
Pitcher2	0.94
Pitcher3	0.67
Pitcher4	0.87
Pitcher5	0.88

## K-Nearest Neighbor

The same six separate model approach was used to train and test the data using K-NN. The K-NN algorithm greatly improved the predictive performance for the total model and each of the separate pitcher models (other than Pitcher 3). Total model accurately predicted 91% of the pitch type in the testing set with Pitcher 1 (96%), Pitcher 2 (96%), Pitcher 4 (93%), and Pitcher 5 (90%) all having greater accuracy than the decision tree model performance.

Table 7: K-NN Model Performance

Model	Accuracy
Total Model	0.91
Pitcher1	0.96
Pitcher2	0.95
Pitcher3	0.67
Pitcher4	0.93
Pitcher5	0.90

### Support Vector Machine (SVM)

As a result of K-NN resulting in an accuracy score above 90% for each separate pitcher model, a multiclass support vector algorithm was ran on the total model to improve the models performance for predicting Pitcher 3 and Pitcher 6 in the testing set. The SVM resulted in a slight improvement in overall model performance (92%) compared to the K-NN total model.

### Final Model Results and Predictions

**Step 4:** The separate K-NN models for Pitcher 1, Pitcher 2, Pitcher 4, and Pitcher 5 reported accuracy scores above 90% (Table 7). Therefore, it was decided to use the total training data for each of the four pitchers to train K-NN models and make final pitch type predictions for these four pitchers in the test data set.

SVM reported the highest predictive accuracy for the total model (92%). It was therefore decided to train SVM on the total training data to make final pitch type prediction for Pitcher 6 as well as Pitcher 3 (due to low observation of training data) in the test data set.

Table 8: Comparing Model Performance

	Total.Model	Pitcher.1	Pitcher.2	Pitcher.3	Pitcher.4	Pitcher.5
Decision Tree Model	0.84	0.93	0.94	0.67	0.87	0.88
K-NN Model	0.91	0.96	0.95	0.67	0.93	0.90
SVM Model	0.92	NA	NA	NA	NA	NA

Table 9: Predicted Model Decision

Variables	Description
Pitcher 1	K-NN: Pitcher specific model
Pitcher 2	K-NN: Pitcher specific model
Pitcher 3	SVM: Total model
Pitcher 4	K-NN: Pitcher specific model
Pitcher 5	K-NN: Pitcher specific model
Pitcher 6	SVM: Total model

**Step 5:** After training K-NN on the total training data for each pitcher. Final predictions were made using each of the four pitchers separate K-NN models. SVM was trained on the total training data and final predictions were made for Pitcher 3 and Pitcher 6. When final predictions were made for each pitcher, the data was merged together to produce a final data set of all pitcher's with their predicted pitcher type.

**Step 6:** The predicted results were displayed along with the actual (i.e., training) data by pitch type and pitcher to visualize if patterns match. Although it appears that velocity had decreased from years 1-2 to year

3 (91, 92 mph vs 87, 89mpg) overall patterns appears similar (e.g., pitch 7 had the overall lowest spin rate, pitch 2 the largest vertical break). Interestingly, it appears that Pitcher 3 and Pitcher 6 are both left-handed pitchers due to both having an initial release point on the first base side of the rubber. This may reduce accuracy rating due to the fact that the data was essentially training on right-handed pitchers to predict pitch type for a left-handed pitcher.

Table 10: Years 1-2

type	mph	spin	breakx	breakz	initx	initz	ext
2	76.99	2512.84	4.89	-6.84	-1.77	5.95	6.19
3	82.46	1867.16	1.06	-0.03	-1.38	5.75	6.21
4	83.82	1364.29	-5.61	3.45	-1.74	5.89	6.20
7	84.63	988.92	-2.91	2.48	-1.02	5.99	6.21
8	88.90	2346.13	1.23	4.71	-1.81	5.81	6.21
9	91.15	2065.17	-7.13	6.91	-1.83	5.86	6.20
10	92.14	2131.53	-3.19	9.45	-1.63	5.94	6.19

Table 11: Final Predictions

PredictedPitchType	mph	spin	breakx	breakz	initx	initz	ext
2	75.67	2686.30	5.32	-5.93	-1.94	6.04	6.20
3	83.63	2008.46	4.93	4.61	3.68	6.41	6.22
4	82.13	1476.81	-6.11	2.43	-1.93	6.03	6.21
7	84.31	1050.90	-1.44	2.47	-0.60	6.03	6.17
8	86.48	2274.63	3.97	6.01	1.63	6.22	6.21
9	87.52	2125.88	-3.19	5.37	-0.43	6.03	6.20
10	89.14	2223.90	-2.05	8.91	-1.06	6.01	6.20



Table 12: Actual Individual Pitcher by Pitch Type: Years 1-2

Pitcher	type	mph	spin	breakx	breakz	initx	initz	ext
Pitcher1	2	77.18	2951.31	5.83	-6.47	-1.85	6.40	6.18
	4	81.26	1432.34	-6.69	0.90	-1.84	6.41	6.20
	9	88.11	2206.98	-8.49	3.43	-2.06	6.28	6.19
	10	89.38	2232.94	-5.99	7.13	-1.89	6.41	6.19
Pitcher2	2	79.73	2574.14	5.55	-6.77	-2.18	5.86	6.20
	4	87.64	1619.41	-5.49	3.16	-2.35	5.74	6.19
	9	93.99	2208.27	-6.27	7.54	-2.26	5.76	6.21
	10	93.99	2241.33	-1.67	8.99	-2.27	5.86	6.20
Pitcher3	3	84.72	2044.59	-0.10	4.39	3.89	6.66	6.21
	9	86.87	2041.95	9.51	4.94	4.14	6.44	6.22
	10	87.67	2098.65	4.79	8.58	4.07	6.51	6.14
Pitcher4	2	81.27	2624.06	4.39	-2.73	-1.92	5.85	6.20
	4	87.02	1430.26	-7.69	3.31	-2.15	5.69	6.18
	8	88.95	2490.01	1.95	4.42	-1.99	5.85	6.21
	9	93.42	2309.79	-7.35	7.70	-2.08	5.86	6.19
	10	93.36	2336.79	-4.69	9.77	-1.97	5.92	6.20
Pitcher5	2	72.03	2167.26	3.96	-9.05	-1.24	5.86	6.19
	3	82.44	1865.42	1.07	-0.08	-1.43	5.75	6.21
	4	82.32	1211.61	-4.57	4.66	-1.33	5.81	6.20
	7	84.63	988.92	-2.91	2.48	-1.02	5.99	6.21
	8	88.81	2068.38	-0.14	5.27	-1.48	5.75	6.20
	9	90.45	1927.53	-7.10	7.43	-1.57	5.78	6.21
	10	90.93	1981.33	-3.10	9.71	-1.17	5.96	6.19

Table 13: Predicted Individual Pitcher by Pitch Type: Year 3

Pitcher	PredictedPitchType	mph	spin	breakx	breakz	initx	initz	ext
Pitcher1	2	76.29	2940.96	5.74	-6.35	-1.85	6.40	6.20
	4	80.56	1432.89	-6.44	0.71	-1.83	6.40	6.23
	9	86.66	2209.15	-7.17	3.24	-2.05	6.29	6.20
	10	87.47	2271.93	-4.56	6.26	-1.89	6.39	6.20
Pitcher2	2	73.48	2573.31	5.55	-6.74	-2.18	5.85	6.19
	4	81.68	1635.80	-5.51	3.40	-2.35	5.74	6.20
	9	87.44	2186.73	-6.68	7.35	-2.25	5.76	6.20
	10	87.67	2237.46	-2.35	8.76	-2.26	5.84	6.21
Pitcher3	2	84.13	2269.58	5.47	4.57	4.21	6.41	6.12
	3	83.79	2024.08	5.33	5.03	4.10	6.46	6.22
	4	85.06	1694.66	10.66	5.92	4.00	6.53	6.21
	7	84.75	1623.93	10.24	5.84	4.01	6.45	6.24
	8	85.39	2135.16	5.38	7.17	4.11	6.49	6.20
	9	83.87	2156.94	2.97	5.60	4.14	6.46	6.17
	10	85.94	2083.01	5.63	8.41	4.06	6.51	6.20
Pitcher4	2	80.47	2620.86	4.38	-2.71	-1.91	5.85	6.21
	4	86.10	1442.49	-7.76	3.63	-2.15	5.70	6.21
	8	88.03	2499.23	2.07	4.32	-1.99	5.85	6.21
	9	92.55	2298.29	-7.29	7.27	-2.04	5.86	6.19
	10	92.57	2340.87	-4.82	9.68	-1.99	5.91	6.19
Pitcher5	2	71.85	2195.99	3.96	-8.94	-1.24	5.87	6.20
	3	81.76	1825.41	0.33	-0.33	-1.31	5.78	6.20
	4	82.04	1225.32	-4.70	4.70	-1.33	5.79	6.17
	7	84.27	1002.53	-2.42	2.19	-0.98	6.00	6.17
	8	88.30	2055.71	-0.50	4.75	-1.50	5.70	6.22
	9	90.07	1943.43	-7.33	7.49	-1.60	5.76	6.22
	10	90.53	1967.73	-3.23	9.61	-1.16	5.96	6.20
Pitcher6	9	87.03	2021.60	3.21	5.51	2.09	5.97	6.20
	10	91.60	2048.80	3.13	10.87	2.17	6.10	6.10