

# STA 531: Final Project

## Predicting Pitches in Baseball

Sarah Normoyle, Drew Jordan, Gonzalo Bustos

May 5, 2016

## 1 Introduction

Sabermetrics, or the analysis of baseball, has been widely used by statisticians as a way to analyze the performance of baseball players and baseball teams. The main focus of sabermetrics has been to compare the performance of individual players. There are also other areas of baseball that would be also for players and coaches. The motivation behind this project is that it would be greatly valuable for a batter to be able to know what pitch to expect before a pitch is thrown. This project attempts to be able to predict the following pitch given a sequence of pitches by a pitcher and a set of covariates about the current game play.

## 2 Data

This data is found publicly on [mlb.com](http://mlb.com), and we obtained it from a hackathon in Baltimore for the baseball team the Orioles. The data set contains pitching data on all MLB teams over the course of 3 years, from 2013 to 2015. There are 2,114,497 observations of pitches in the data set, and 18 different variables.

## 3 Methods

This project is aimed at applying statistical methods to analyze the pitching sequence of MLB pitchers. Various methods were implemented, and the results were compared to each other. The two main statistical techniques that were used were Markov Models and Multinomial Logistic Regression to predict the sequence of pitches. The sections below go further into detail into the methods used. To be able to compare the different methods, we used cross-validation for each method and used percent accuracy in the predictions as the comparison metric. For each of these methods, we focused on applying the techniques to one particular pitcher.

### 3.1 Sampling from Probability Vectors

The first naive method that was used as a baseline was obtaining the sample probability vectors of throwing the various pitches. This method does not consider any of the covariates. After obtaining a sample probability vector of throwing each pitch, which is obtained by getting the counts for each pitch and dividing by the total number of pitches, we can then sample from the different pitch possibilities with probability of the probability vector. Once these predicted pitches are calculated, we obtain the percent accuracy of predictions. We ran this 100 times as there is inherent randomness in the sampling.

### 3.2 Markov Model

Next, a simple Markov model was used. This Markov model only used an initial probability vector and a transition matrix to predict the next pitch in a sequence. The transition matrix was first created from a

known set of pitches, and then it can be used to predict the next pitch in a sequence. Cross validation was also applied with this method by training, or creating the transition matrix, on a portion of the data, and then testing this Markov model on the held out set. The percent accuracy in predictions was once again used as the comparison metric.

### 3.3 Hidden Markov Model

After applying a simple Markov model, a more complex hidden Markov model was implemented to the pitching sequences. Given a sequence of pitches, to obtain the optimal parameters for the model, the Baum-Welch algorithm was used. Once these parameters were determined, to obtain the probabilities of the following observation given the previous observations,  $p(x_{n+1}|x_{1:n})$ , the Forward Algorithm was then implemented. Both algorithms are described more in detail below.

#### 3.3.1 Baum-Welch

The Baum-Welch algorithm applies expectation maximization to hidden Markov models to obtain parameters for the Hidden Markov Model, which are the initial probability vector, the transition matrix from state to state, and the emission matrix. The Baum-Welch algorithm is an iterative algorithm that uses the forward-backward algorithm at each iteration to estimate these parameters.

The forward-backward algorithm is as follows:

In the forward algorithm, we sum over  $z_1, z_2, \dots, z_n$  in that order, and derive a recursion for computing  $p(x_{1:j}, z_j)$  for each  $z_j = 1, \dots, m$  and each  $j = 1, \dots, n$ . In the backward algorithm, we sum over  $z_n, z_{n-1}, \dots, z_1$ , and derive a recursion for computing  $p(x_{j+1:n}|z_j)$  for each  $z_j = 1, \dots, m$  and each  $j = 1, \dots, n$ .

The forward-backward algorithm was implemented as shown below. The log of the probabilities were taken in order to deal with arithmetic underflow/overflow and R's inability to store such a low probability. In addition, the log-sum-exp trick was used to deal with a similar problem.

Forward algorithm:

1. For each  $z_1, \dots, m$ , compute  $g_1(z_1) = \log p(z_1) + \log p(x_1|z_1)$
2. For each  $j = 2, \dots, n$  for each  $z_j = 1, \dots, m$ , compute:

$$\log s_j(z_j) = g_j(z_j) = \log \sum_{z_{j-1}} \exp[g_{j-1}(z_{j-1}) + \log p(z_j|z_{j-1}) + \log p(x_j|z_j)]$$

3.  $\log p(x_{1:n}) = \log \sum_{z_n} \exp(g_n(z_n))$

And  $g_j(z_j) = \log p(x_{1:j}, z_j)$

Backward algorithm:

1. For each  $z_n = 1, \dots, m$ , define  $r_n(z_n) = 0$
2. For each  $j = n - 1, n - 2, \dots, 1$ , for each  $z_j = 1, \dots, m$  compute:

$$r_j(z_j) = \log \sum_{z_{j+1}} \exp(\log p(z_{j+1}|z_j) + \log p(x_{j+1}|z_{j+1}) + r_{j+1}(z_{j+1}))$$

And  $r_j(z_j) = \log p(x_{j+1:n}|z_j)$

The Baum-Welch algorithm is implemented as follows:

Using the following formula:

$$\begin{aligned}\gamma_{ti} &= P(Z_t = i|x) \\ \beta_{tij} &= P(Z_{t-1} = i, Z_t = j|x) \\ \pi_j &= \frac{\gamma_{1i}}{\sum_{j=1}^m \gamma_{1j}} \\ T_{ij} &= \frac{\sum_{t=2}^n \beta_{tij}}{\sum_{t=1}^{n-1} \gamma_{ti}}\end{aligned}$$

The algorithm is:

1. Randomly initialize  $\pi, T$ , and  $\phi = (\phi_1, \dots, \phi_m)$
2. Iteratively repeat the following two steps, until convergence:
  - (a) E-step: Compute the  $\gamma$  and  $\beta$  using the forward-backward algorithm.
  - (b) M-step: Update  $\pi, T$ , and  $\phi$  using the formulas above.

This algorithm was implemented in Python, and the results from the Baum-Welch were used in the Forward algorithm to do probabilistic inference, which is described below.

### 3.3.2 Forward Algorithm

Once the parameters are estimated for the transition matrix, the emission matrix, and the initial probability matrix, we can combine these parameters with results from the forward algorithm to get probabilities for the next observation given the previous observation.

We can predict  $x_{j+1}$  and  $x_{1:j}$  using:

$$\begin{aligned}p(x_{j+1}|x_{1:j}) &\propto p(x_{1:j}, x_{j+1}) = \sum_{z_j, z_{j+1}} p(x_{1:j}, x_{j+1}, z_j, z_{j+1}) \\ &= \sum_{z_j, z_{j+1}} p(x_{1:j}, z_j) p(z_{j+1}|z_j) p(x_{j+1}|z_{j+1})\end{aligned}$$

Cross validation was also applied in this setting. The Baum-Welch algorithm was applied to the training set. Then as we ran through a sequence of observations in the testing set, we can use the forward algorithm and the estimate parameters to get probabilities of the next observation. These predictions are compared to the true predictions in order to obtain an estimated percent accuracy in predictions.

## 3.4 Multinomial Logistic Regression

After Markov models were implemented, multinomial logistic regression models were then fit to predict the pitch given a set of covariates. Multinomial logistic regression is a classification scheme that generalizes logistic regression to a multivariate scheme. Because a pitcher can throw multiple types of pitches at a given observation, we have a vector of possible pitches.

Given the multinomial data with  $J$  categories and the  $p$ -dimensional predictor variables, we can forecast in which  $j$  category a future data point  $y^*$  at the predictor  $x^*$  will be.

Multinomial logistic regression can be understood as a set of independent binary regressions. If we have  $J$  possible outcomes, we can imagine running  $J - 1$  binary regression models, which are compared against the one pivot outcome.

$$\begin{aligned}\frac{\Pr(Y_i = 1)}{\Pr(Y_i = K)} &= \beta_1 \cdot X_i \\ \frac{\Pr(Y_i = 2)}{\Pr(Y_i = K)} &= \beta_2 \cdot X_i \\ &\dots\dots\dots \\ \frac{\Pr(Y_i = K - 1)}{\Pr(Y_i = K)} &= \beta_{K-1} \cdot X_i\end{aligned}$$

For each possible outcome, there are separate vectors of regression coefficients. Therefore, we can calculate the probability of observing a category  $j$  after at a time occurrence as:

$$P(y^* = j | x^*, \beta, n^* = 1) = e^{x^* \beta_j} / \sum_{k=1}^J e^{x^* \beta_k}.$$

Once we calculate these probabilities for a set of covariates, we can predict what the pitch will be for the set of covariates.

For this dataset, there were specific variables that were intended to be used in the model, as they have known to make a difference in pitching. There were also some attempts at variable selection by choosing variables that created the highest percent accuracy when implemented on the testing set.

### 3.5 Cross Validation

Cross Validation was implemented for the various methods. Cross Validation

## 4 Results

blah

### 4.1 Exploratory Data Analysis

blah

### 4.2 Predicting and Cross Validation

The various methods descried in the Methods section were implemented for the pitcher Clayton Kershaw. He only pitches four different type of pitches, and he has numerous pitch observations in the dataset. The total percent accuracy was calculated for each method.

Results from testing within sample:

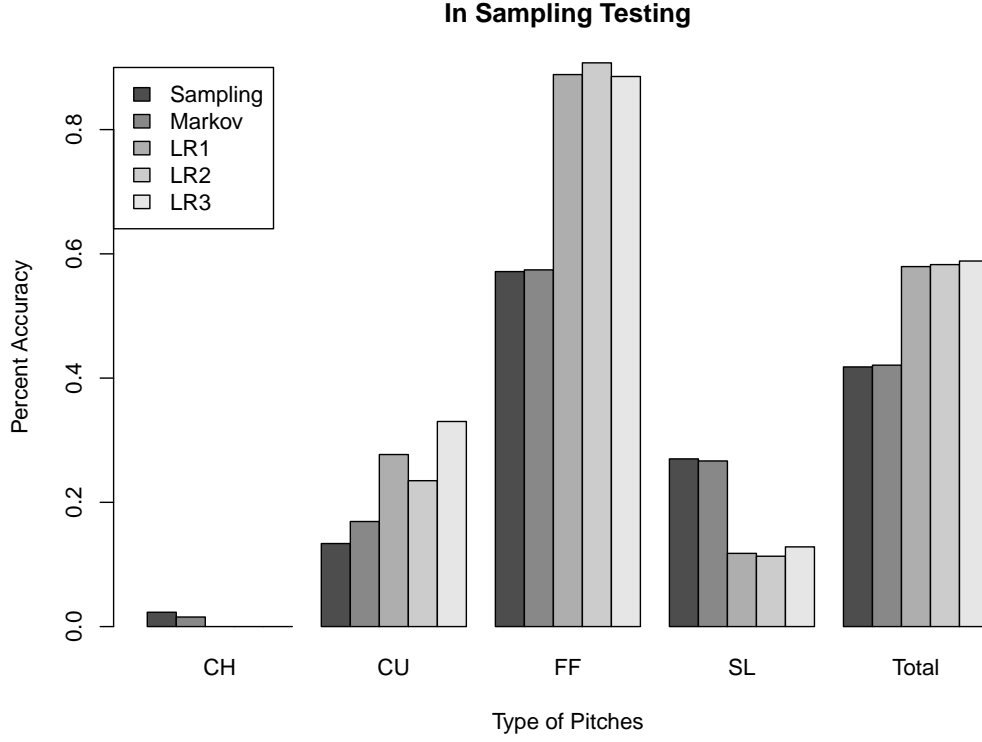
	CH	CU	FF	SL	Total
Sampling	0.0231	0.1337	0.5715	0.2700	0.4180
Markov	0.0154	0.1691	0.5742	0.2666	0.4208
LR1	0.0000	0.2770	0.8886	0.1178	0.5795
LR2	0.0000	0.2349	0.9074	0.1132	0.5827
LR3	0.0000	0.3301	0.8855	0.1283	0.5885

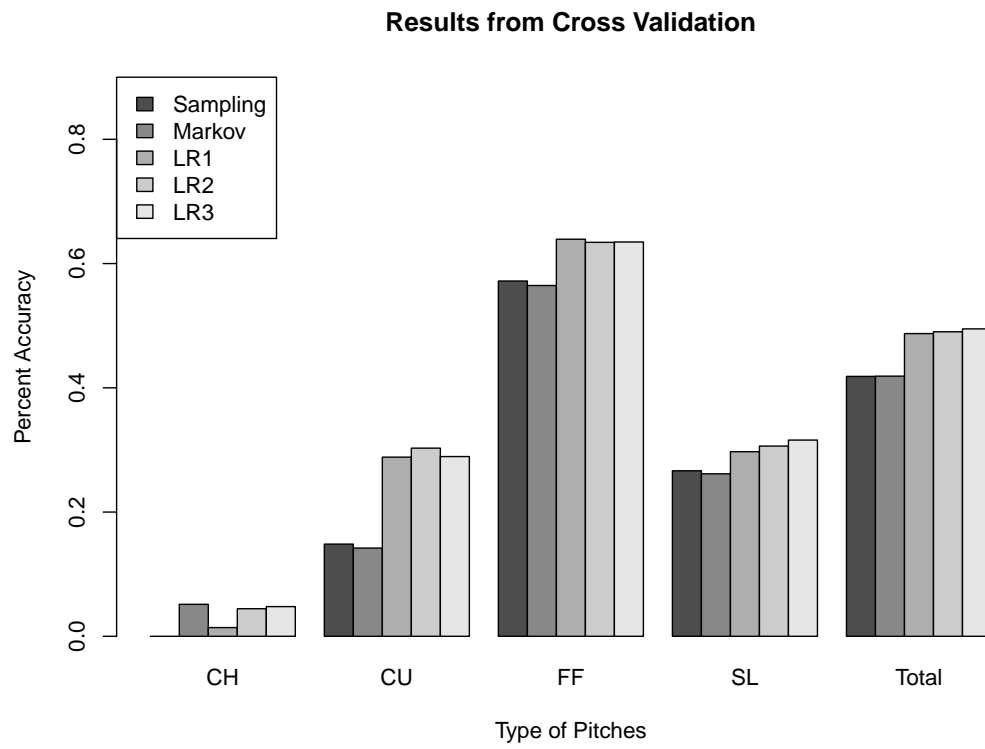
Table 1: Percent Accuracies from In-Sample Testing

Results from Cross Validation:

	CH	CU	FF	SL	Total
Sampling	0.0000	0.1484	0.5719	0.2665	0.4183
Markov	0.0515	0.1420	0.5647	0.2616	0.4188
LR1	0.0140	0.2884	0.6392	0.2972	0.4873
LR2	0.0444	0.3029	0.6342	0.3062	0.4902
LR3	0.0478	0.2893	0.6347	0.3159	0.4949

Table 2: Percent Accuracies from Cross Validation





## 5 Conclusion

blah