

A Bayesian Approach to Fantasy Baseball

Adam Mandelson
ISYE6240
25-April-2021

Introduction

This project explores weekly matchup data for the fantasy baseball league for the 2016-2019 seasons. Parts of this analysis are motivated by [Pedar Coyle's](#) [Rugby Analysis](#) case study, and [PyMC3's](#) [GLM Logistic Regression](#) example. The data was collected using the [Yahoo Fantasy Baseball](#) API and is available in the `full_data.csv` file. The league consists of 14 teams and each season usually consists of 21 weekly matchups. Scoring is based on 18 statistical categories. Weekly "scores" are based on the head-to-head matchup in each category. Category wins, losses and ties are compiled each week and throughout the season to determine the league standings. The 2016-2019 seasons were chosen because the categories were consistent throughout those four seasons. They consist of the following counting (aggregate) categories as identified in the data:

Batting: `b_BB`: Walks; `b_HR`: Home Runs; `b_K`: Strikeouts; `b_NSB`: Net Stolen Bases; `b_R`: Runs; `b_RBI`: Runs Batted In.
Pitching: `p_BB`: Walks; `p_IP`: Innings Pitched; `p_K`: Strikeouts; `p_L`: Losses; `p_NSVH`: Net Saves and Holds; `p_OS`: Quality Starts; `p_W`: Wins.

There are rate categories:

Batting: `b_AVG`: Batting Average; `b_OBP`: On-Base Percentage; `b_SLG`: Slugging Percentage.
Pitching: `p_ERAR`: Earned-run Average; `p_WHIP`: Walks plus Hits per Innings Pitched.

As noted above, the season standings are not influenced by the overall weekly winning, and consist of an aggregate score. For example, in a weekly matchup of two teams, the result is not purely a win or a loss. The "winning" team might win 10 categories, lose six and tie two, and the "score" of 10-6-2 is simply added to their standings. Despite this, Yahoo's API includes a "weekly winning team" category, and for the sake of exploring the data and Bayesian approaches, the concept of a team winning the weekly matchup will be used to develop exploratory models.

Motivation

Ideally, this analysis could be used to infer a parameter or a threshold that would lead to more wins in each category. One thought was to identify key categories to focus on, and develop a predictive model for weekly matchup wins based on individual categories. However, I found that the variances in the data and the large number of categories meant that each category has a relatively small individual impact on the overall weekly likelihood of a "win".

Likewise, as we'll see below, the data is highly-correlated in ways that make sense: a hit results in an increase in batting average, on-base percentage, and slugging percentage; a home run results in an increase in home runs, runs, runs batted in, slugging percentage, and on-base percentage.

As such, this analysis aims to explore the data and the Bayesian logistic regression toolset in the `PyMC3` library, and to make general observations about the data and how the models could be used or fine-tuned.

Lastly, I should note that the high-dimensionality in the data meant that attempts to elicit priors were ineffective, and sometimes caused errors in the computations. Based on experiences this semester in OpenBUGs in which we defined uninformative priors, I felt it was a good idea to defer to `PyMC3's` defaults for [generalized linear models](#), which are a flat, uninformative prior distribution for the intercept, and a Normal distribution with mean=0 and precision=1.0E-6 for all regressors.

Data Import and Analysis

```
In [1]: import pandas as pd
import os
import warnings
warnings.filterwarnings('ignore')

# Import data directories, or load file if in the same directory
full_data_csv = 'full_data.csv'

from config import DATA_DIRECTORY, LEAGUE_DATA_DIRECTORY
league_directory = LEAGUE_DATA_DIRECTORY
data_directory = DATA_DIRECTORY
df = pd.read_csv(os.path.join(data_directory, "full_data.csv"))
df.info()
except FileNotFoundError:
    df = pd.read_csv("../full_data.csv")
    df.info()

# Matplotlib inline
%matplotlib inline

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21643 entries, 0 to 21642
Data columns (total 8 columns):
#   Column      Non-Null Count  Dtype
---  --
0   season      22644 non-null    int64
1   week        22644 non-null    int64
2   n_days      22644 non-null    int64
3   team_id     22644 non-null    int64
4   category    22644 non-null    object
5   value       22644 non-null    float64
6   won_category 22644 non-null    float64
7   won_week    22644 non-null    bool
dtypes: bool(1), float64(2), int64(4), object(1)
memory usage: 1.2+ MB

The data exported from Yahoo is long, 22644 lines, each representing a single category value for the season, week and team_id. The 'won_category' field is an integer value of 1 or 0 depending on whether the team won (1), lost (0) or tied (0).

A bit of adjusting is necessary. The 'won_category' field for each line has NA values that will be replaced with zeroes.

The season, week and team_id fields will be changed to strings to create an 'id' field that identifies all of the categories and values for a specific team during a specific week within a specific season. This will help pivot the data to look at the overall weekly winners.
```

```
In [2]: # Fill in NA values.
df['won_category'] = df['won_category'].fillna(0) # There are ties

# Create a unique 'id' to turn long data for each category for each week for each team into weekly team data
df['season'] = df['season'].astype(str)
df['week'] = df['week'].astype(str)
df['team_id'] = df['team_id'].astype(str)
df['id'] = df['season'].str.cat(df.loc[:, ['week', 'team_id']], sep='_') # Ex: 2019_5_14 for 19 categories

# Change the won category field back to bool.
df['won_category'] = df['won_category'].astype(bool)

# Example:
df.sample(5)
```

```
Out[2]:
```

	season	week	n_days	team_id	category	value	won_category	won_week	id
20592	2019	15	7	9	b_R	55.000	False	True	2019_15_9
12586	2018	6	7	10	b_K	59.000	False	False	2018_6_10
10375	2017	20	7	4	b_SLG	0.383	False	False	2017_20_4
18553	2019	7	7	4	p_K	79.000	True	False	2019_7_4
18886	2019	8	7	12	b_K	61.000	False	False	2019_8_12

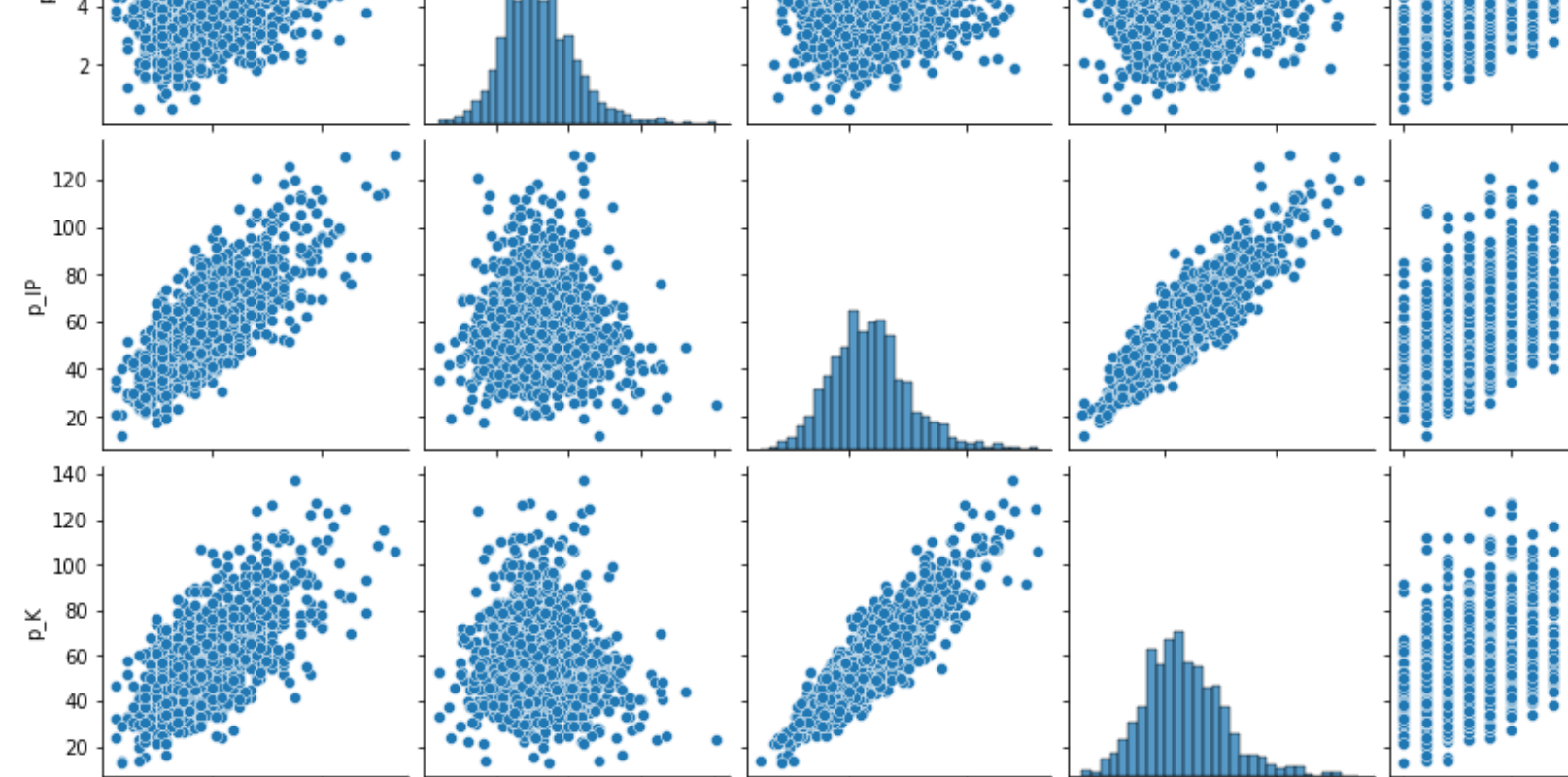
```
In [3]: # Pivot the data
df_pivoted = df.pivot(index='id', columns='category', values='value')
df_pivoted['won_week'] = df.groupby('id')['won_week'].mean()
```

```
In [4]: import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

# A correlation matrix
corr = df_pivoted.corr()
```

```
mask = np.zeros_like(corr, dtype=np.bool)
mask[np.triu_indices_from(mask)] = True
i, ax = plt.subplots(figsize=(8, 6))
cmmap = sns.diverging_palette(10, 220, as_cmap=True)
sns.heatmap(corr, mask=mask, cmap=cmmap, vmax=0.3, linewidth=0.5, cbar_kws={"shrink": 0.5}, ax=ax)
```

```
Out[4]: <AxesSubplot: xlabel='category', ylabel='category'>
```



As discussed earlier, there are some clear correlations here, though the highest is about 0.35 (`b_OBP`). `b_AVG` and `b_SLG` are logically correlated, as well.

```
In [5]: np.absolute(corr['won_week']).sort_values(ascending=False)
```

```
Out[5]:
```

category	won_week
b_OBP	1.000000
b_AVG	0.351225
b_SLG	0.314886
b_R	0.306156
b_RBI	0.268607
b_K	0.241708
p_WHIP	0.238372
p_ERAR	0.230599
p_OS	0.216396
b_R	0.209132
p_W	0.199857
b_BB	0.183897
b_NSB	0.152835
p_K	0.144011
p_IP	0.136806
p_L	0.093465
b_K	0.071859
b_NSVH	0.058432
b_BB	0.048027
Name: won_week, dtype: float64	

```
In [6]: # Create some additional datasets for further exploration
df_batting = df_pivoted.loc[:, [c for c in df_pivoted.columns if 'b_' in c]]
df_pitching = df_pivoted.loc[:, [c for c in df_pivoted.columns if 'p_' in c]]

counts = df_pivoted.loc[:, [c for c in df_pivoted.columns if not 'z' in c]]
rates = df_pivoted.loc[:, [c for c in df_pivoted.columns if 'z' in c]]
batting_counts = counts.loc[:, [c for c in counts.columns if 'b_' in c]]
pitching_counts = counts.loc[:, [c for c in counts.columns if 'p_' in c]]

df_batting_counts = df_batting.loc[:, [c for c in df_batting.columns if not 'z' in c]]
df_batting_rates = df_batting.loc[:, [c for c in df_batting.columns if 'z' in c]]
df_pitching_counts = df_pitching.loc[:, [c for c in df_pitching.columns if not 'z' in c]]
df_pitching_rates = df_pitching.loc[:, [c for c in df_pitching.columns if 'z' in c]]
```

```
In [7]: # sns.pairplot(df_batting) --> Too large for printing
```

```
In [8]: # A few pitching categories for further exploration.
sns.pairplot(df_pitching.loc[:, :5])
```

```
Out[8]: <seaborn.axisgrid.PairGrid at 0x2030be4bb50>
```



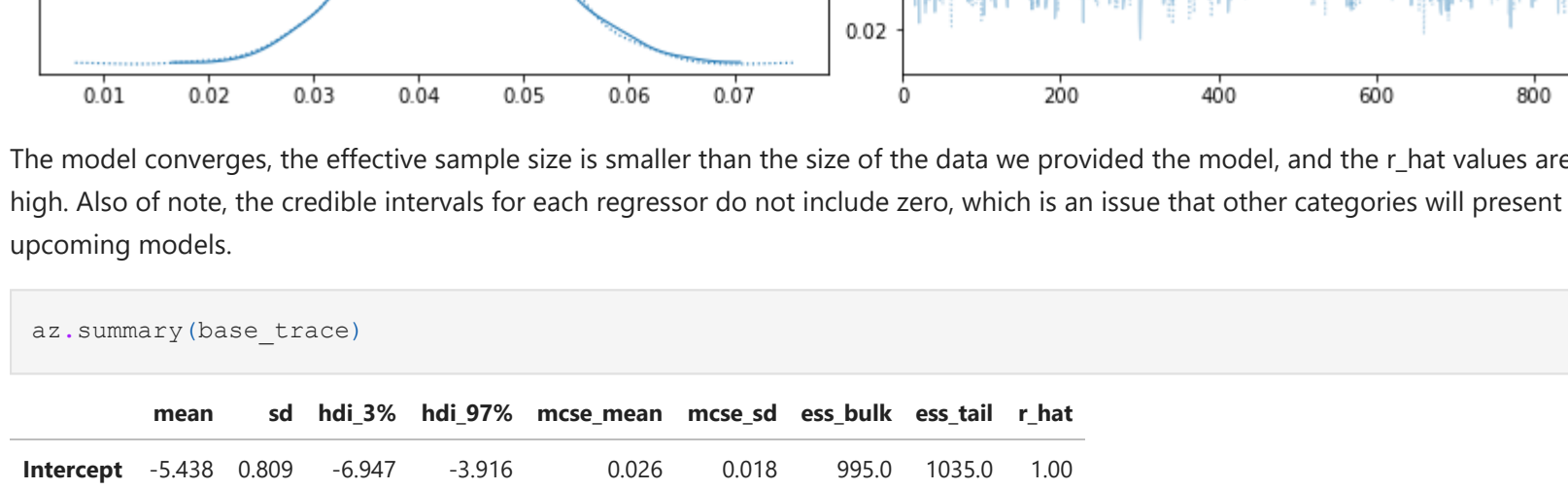
Not surprisingly, the rate categories broadly correlate with Innings Pitched. The losses category has some interesting points in it, as well, as there is a clear ERA threshold as losses increase.

The strong correlation between Innings Pitched and strikeouts is logical, but what's notable is the lesser correlation with walks. It's also worth noting that a few of these categories may be multimodal.

```
In [9]: # sns.pairplot(rates)
```

```
In [10]: # The strongest category-won_week correlation
sns.distplot(df_pivoted['b_OBP'])
```

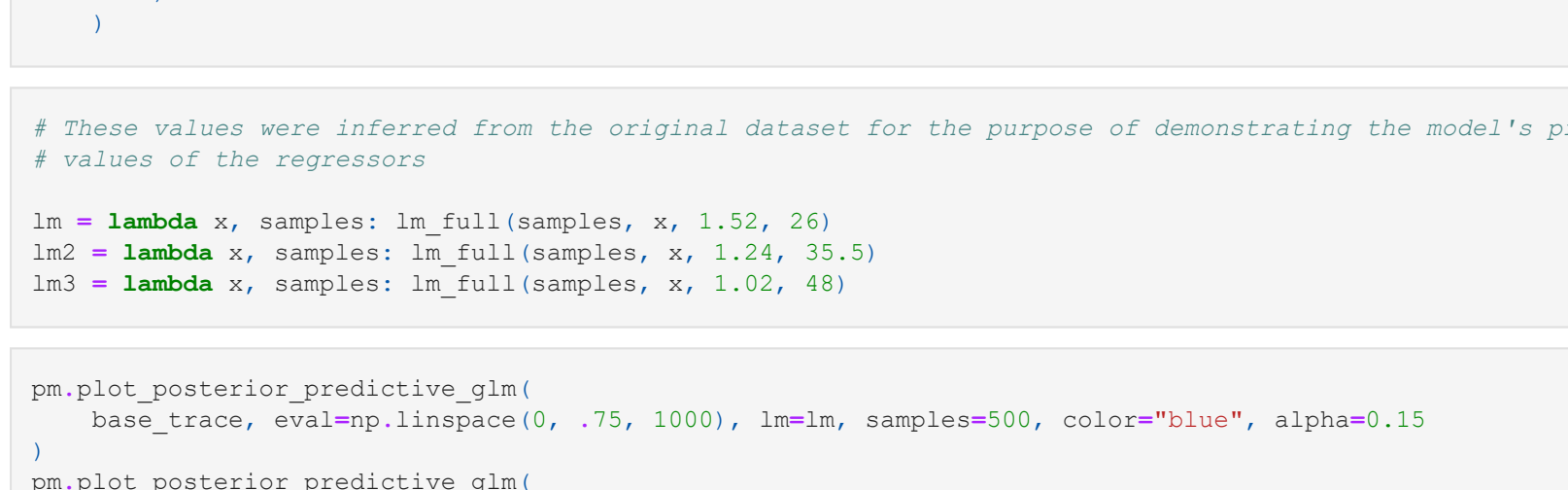
```
Out[10]: <AxesSubplot: xlabel='b_OBP', ylabel='Density'>
```



The distribution appears mostly normal, but may be well-represented by a Beta distribution, as well.

```
In [11]: # The strongest pitch-rate category-won_week correlation
sns.distplot(df_pivoted['p_WHIP'])
```

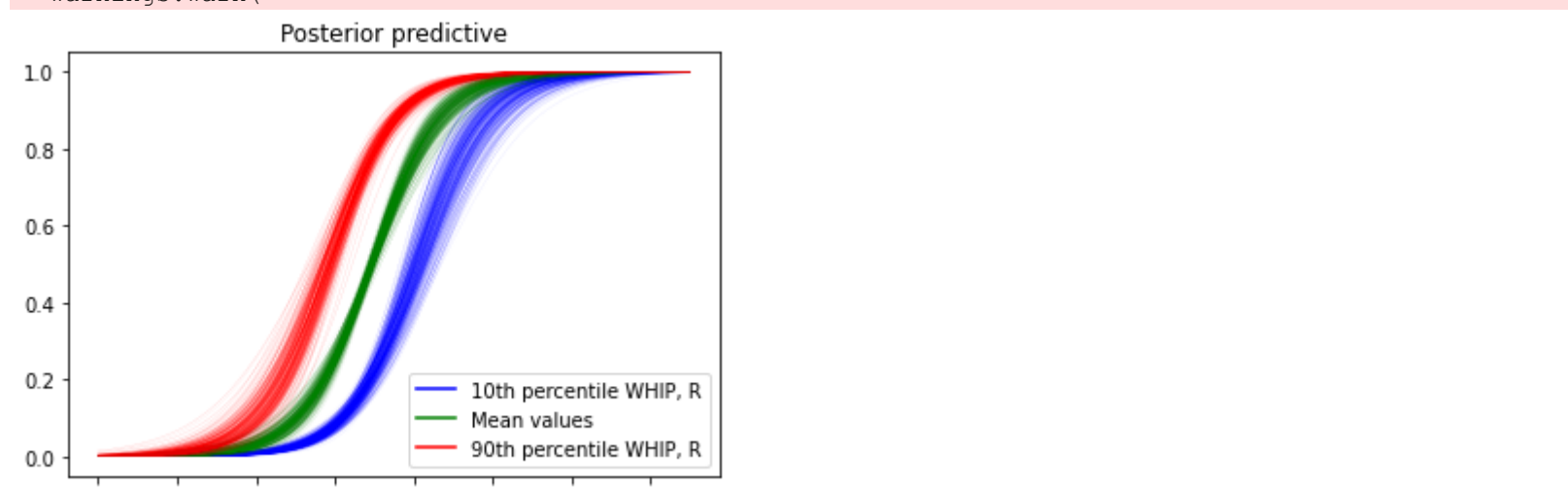
```
Out[11]: <AxesSubplot: xlabel='p_WHIP', ylabel='Density'>
```



`p_WHIP` appears to be normally distributed, but may be multimodal.

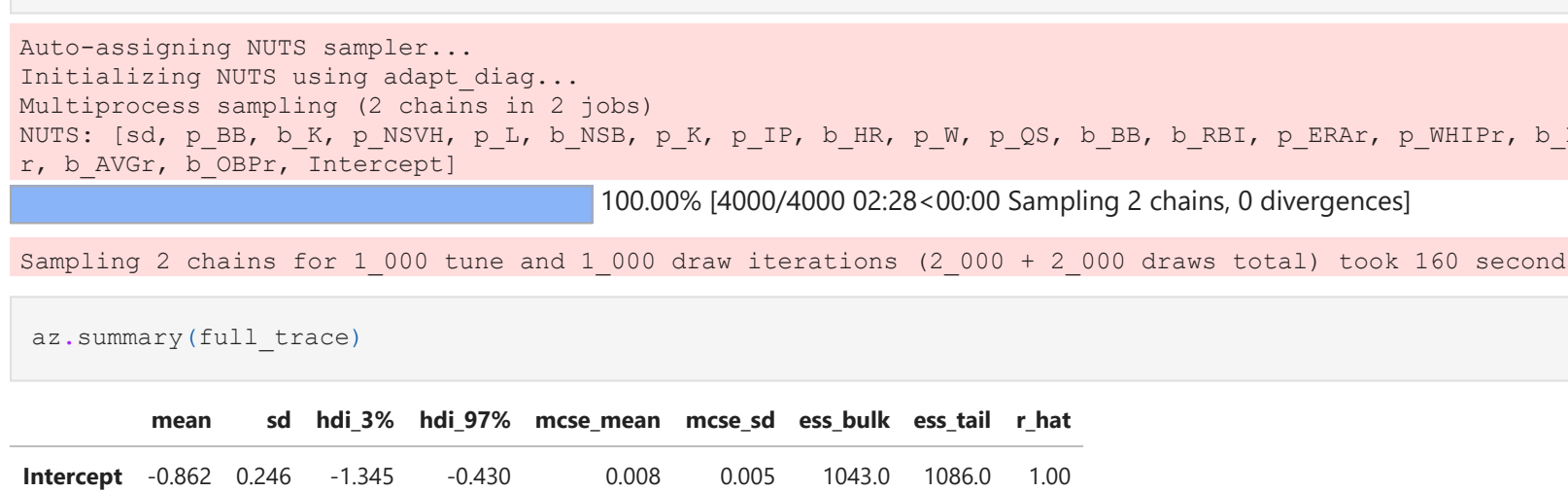
```
In [13]: sns.stripplot(x='won_week', y='b_OBP', data=df_pivoted, jitter=True)
```

```
Out[13]: <AxesSubplot: xlabel='won_week', ylabel='b_OBP'>
```



```
In [14]: sns.stripplot(x='won_week', y='p_WHIP', data=df_pivoted, jitter=True)
```

```
Out[14]: <AxesSubplot: xlabel='won_week', ylabel='p_WHIP'>
```



Again, because of the large number of categories, we're seeing that there isn't a lot we can infer from these categories about how they contribute to weekly wins, though `p_WHIP` does show that, like `ERA`, there appears to be a bit of a threshold where a WHIP above 1.75 or so rarely, if ever, results in a win. This may be a good target for our regression models.

The models

The first model will look at just the three categories we've just explored.

```
In [16]: import pymc3 as pm
import arviz as az

# This models won_week as the likelihood product of bernoulli n trials
# By default the priors on the regressors are very vague at N(0, 1.0E-16)
```

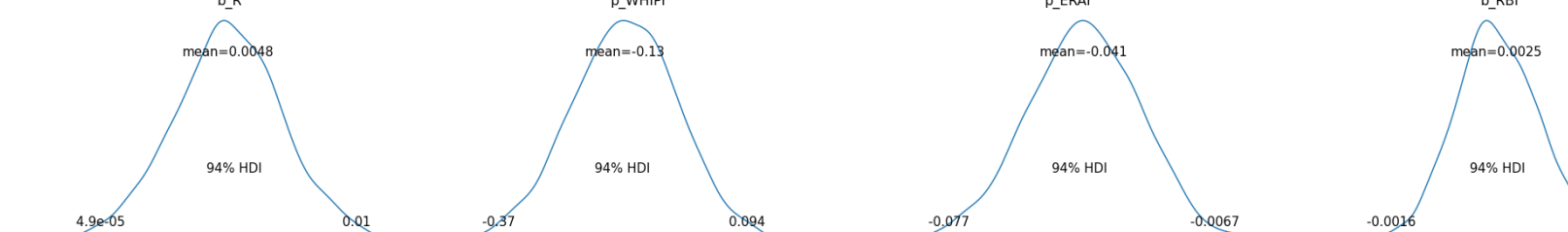
```
with pm.Model() as base_model:
    pm.glm.GLM.from_formula(
        'won_week ~ b_OBP + p_WHIP + b_R', df_pivoted, family=pm.glm.families.Binomial()
    )
    base_trace = pm.sample(init='adapt_diag')
```

Auto-assigning NUTS sampler...
Initializing NUTS using adapt_diag...
Multi-process sampling (2 chains in 2 jobs)
NUTS: [b_R, p_WHIP, b_OBP, Intercept] 100.00% [4000/4000 01:02<00:00 Sampling 2 chains, 0 divergences]

Sampling 2 chains for 1,000 tune and 1,000 draw iterations (2,000 + 2,000 draws total) took 72 seconds.

```
In [17]: az.plot_trace(base_trace, compact=True)
```

```
Out[17]: array([[<AxesSubplot: title='(center: 'Intercept')',
<AxesSubplot: title='(center: 'Intercept')',
<AxesSubplot: title='(center: 'b_OBP')',
<AxesSubplot: title='(center: 'b_OBP')',
<AxesSubplot: title='(center: 'p_WHIP')',
<AxesSubplot: title='(center: 'p_WHIP')',
<AxesSubplot: title='(center: 'b_R')',
<AxesSubplot: title='(center: 'b_R')',
<AxesSubplot: title='(center: 'b_R')',
<AxesSubplot: title='(center: 'b_R')']], dtype=object)
```



The model converges, the effective sample size is smaller than the size of the data we provided the model, and the `r_hat` values are not too high. Also of note, the credible intervals for each regressor do not include zero, which is an issue that other categories will present in upcoming models.

```
In [18]: az.summary(base_trace)
```

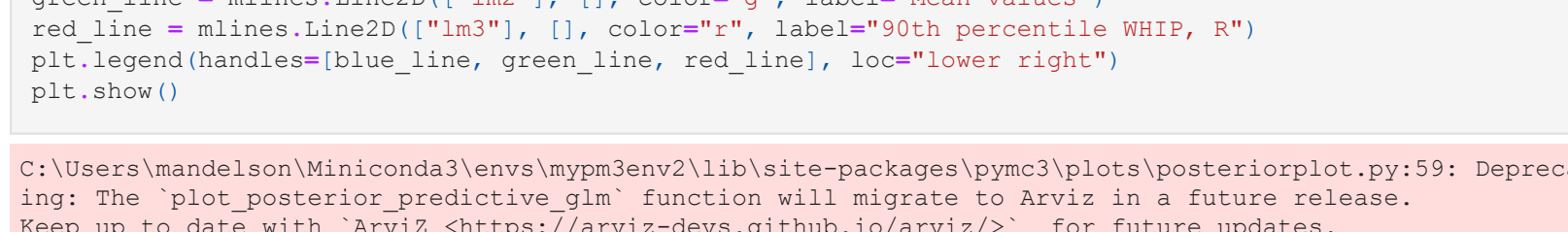
```
Out[18]:
```

	mean	sd	hdi_3%	hdi_97%	mcmc_mean	mcmc_sd	ess_bulk	ess_tail	r_hat
Intercept	-5.438	0.809	-6.947	-3.916	0.026	0.018	995.0	1086.0	1.00
b_OBP	0.474	1.820	1.385	7.443	0.067	0.047	1043.0	1088.0	1.00
b_SLG	-1.262	1.658	-4.350	1.920	0.055	0.039	907.0	1299.0	1.00
b_AVG	0.649	0.643	-0.553	1.822	0.018	0.013	1030.0	1385.0	1.00
b_R	0.005	0.003	0.000	0.010	0.000	0.000	1535.0	1377.0	1.00
p_WHIP	-0.130	0.122	-0.369	0.094	0.004	0.003	1109.0	1219.0	1.00
p_ERAR	-0.041	0.019	-0.077	-0.007	0.001	0.000	1322.0	1377.0	1.00
b_RBI	0.003	0.002	-0.002	0.007	0.000	0.000	1993.0	1659.0	1.00
b_BB	0.001	0.005	-0.008	0.009	0.000	0.000	759.0	1155.0	1.00
p_OS	0.028	0.011	0.009	0.050	0.000	0.000	1674.0	1625.0	1.00
p_W	0.012	0.009	-0.006	0.027	0.000	0.000	1635.0	1415.0	1.00
b_HR	0.003	0.009	-0.013	0.019	0.000	0.000	1385.0	1293.0	1.00
p_IP	0.004	0.002	-0.001	0.008	0.000	0.000	1100.0	1268.0	1.00
p_K	0.000	0.001	-0.003	0.003	0.000	0.000	1829.0	1660.0	1.00
b_NSB	0.024	0.005	0.014	0.033	0.000	0.000	2604.0	1268.0	1.01
p_L	-0.014	0.009	-0.031	0.003	0.000	0.000	1624.0	1414.0	1.00
p_NSVH	0.003	0.005	-0.007	0.013	0.000	0.000	2455.0	1385.0	1.00
b_K	-0.005	0.001	-0.007	-0.002	0.000	0.000	1612.0	1533.0	1.00
p_BB	-0.009	0.003	-0.015	-0.004	0.000	0.000	1417.0	1397.0	1.00
b_K	0.433	0.009	0.416	0.448	0.000	0.000	2035.0	1382.0	1.00

Some of these coefficients have slightly higher `r_hat` values, or low effective sample size values. A number of them contain zero within their confidence interval, though I think with more distinct priors that could be avoided.

```
In [24]: az.plot_posterior(full_trace)
```

```
Out[24]: array([[<AxesSubplot: title='(center: 'Intercept')',
<AxesSubplot: title='(center: 'b_OBP')',
<AxesSubplot: title='(center: 'b_AVG')',
<AxesSubplot: title='(center: 'b_SLG')',
<AxesSubplot: title='(center: 'b_R')',
<AxesSubplot: title='(center: 'p_WHIP')',
<AxesSubplot: title='(center: 'p_ERAR')',
<AxesSubplot: title='(center: 'b_RBI')',
<AxesSubplot: title='(center: 'b_BB')',
<AxesSubplot: title='(center: 'p_OS')',
<AxesSubplot: title='(center: 'p_W')',
<AxesSubplot: title='(center: 'b_HR')',
<AxesSubplot: title='(center: 'p_IP')',
<AxesSubplot: title='(center: 'p_K')',
<AxesSubplot: title='(center: 'b_NSVH')',
<AxesSubplot: title='(center: 'b_K')',
<AxesSubplot: title='(center: 'p_BB')',
<AxesSubplot: title='(center: 'b_R')',
<AxesSubplot: title='(center: 'b_R')']], dtype=object)
```



It's clear from this data that with better values for each category, the model's posterior distribution predicts a higher probability of winning the week. It's also interesting to note that the better values of WHIP and R offset lower OBP values to the point that a weekly win is highly probable right around the mean value for OBP.

Other model runs were less informative, and are not worth including here. The same goes for attempts to model scaled data.

One model that I did want to compare was the full dataset in a logistic regression model:

```
In [22]: with pm.Model() as full_model:
    pm.glm.GLM.from_formula(
        'won_week ~ b_OBP + b_AVG + b_SLG + b_R + p_WHIP + p_ERAR + b_RBI + b_BB + p_OS + p_W + b_HR + p_IP + b_K + p_NSVH + p_L + b_NSB + p_K + p_IP + b_HR + p_W + p_OS + b_BB + b_RBI + p_WHIP + b_R + b_SLG + b_AVG + b_OBP, Intercept)
    full_trace = pm.sample(init='adapt_diag', cores=2)
```

Auto-assigning NUTS sampler...
Initializing NUTS using adapt_diag...
Multi-process sampling (2 chains in 2 jobs)
NUTS: [b_OBP, b_AVG, b_SLG, b_R, p_WHIP, p_ERAR, b_RBI, b_BB, p_OS, p_W, b_HR, p_IP, b_K, p_NSVH, p_L, b_NSB, p_K, p_IP, b_HR, p_W, p_OS, b_BB, b_RBI, p_WHIP, b_R, b_SLG, b_AVG, b_OBP, Intercept] 100.00% [4000/4000 02:28<00:00 Sampling 2 chains, 0 divergences]

Sampling 2 chains for 1,000 tune and 1,000 draw iterations (2,000 + 2,000 draws total) took 160 seconds.

```
In [23]: az.summary(full_trace)
```

```
Out[23]:
```

	mean	sd	hdi_3%	hdi_97%	mcmc_mean	mcmc_sd	ess_bulk	ess_tail	r_hat
Intercept	0.862	0.246	-1.385	-0.430	0.067	0.005	1043.0	1086.0	1.00
b_OBP	0.474	1.820	1.385	7.443	0.067	0.047	1043.0	1088.0	1.00
b_SLG	-1.262	1.658	-4.350	1.920	0.055	0.039	907.0	1299.0	1.00
b_AVG	0.649	0.643	-0.553	1.822	0.018	0.013	1030.0	1385.0	1.00
b_R	0.005	0.003	0.000	0.010	0.000	0.000	1535.0	1377.0	1.00
p_WHIP	-0.130	0.122	-0.369	0.094	0.004	0.003	1109.0	1219.0	1.00
p_ERAR	-0.041	0.019	-0.077	-0.007	0.001	0.000	1322.0	1377.0	1.00
b_RBI	0.003	0.002	-0.002	0.007	0.000	0.000	1993.0	1659.0	1.00
b_BB	0.001	0.005	-0.008	0.009	0.000	0.000	759.0	1155.0	1.00
p_OS	0.028	0.011	0.009	0.050	0.000	0.000	1674.0	1625.0	1.00
p_W	0.012	0.009	-0.006	0.027	0.000	0.000	1635.0	1415.0	1.00
b_HR	0.003	0.009	-0.013	0.019	0.000	0.000	1385.0	1293.0	1.00
p_IP	0.004	0.002	-0.001	0.008	0.000	0.000	1100.0	1268.0	1.00
p_K	0.000	0.001	-0.003	0.003	0.000	0.000	1829.0	1660.0	1.00
b_NSB	0.024	0.005	0.014	0.033	0.000	0.000	2604.0	1268.0	1.01
p_L	-0.014	0.009	-0.031	0.003	0.000	0.000	1624.0	1414.0	1.00
p_NSVH	0.003	0.005	-0.007	0.013	0.000	0.000	2455.0	1385.0	1.00
b_K	-0.005	0.001	-0.007	-0.002	0.000	0.000	1612.0	1533.0	1.00
p_BB	-0.009	0.003	-0.015	-0.004	0.000	0.000	1417.0	1397.0	1.00
b_K	0.433	0.009	0.416	0.448	0.000	0.000	2035.0	1382.0	1.00

Some of these coefficients have slightly higher `r_hat`

