

# Predicting Top 10 HR Hitters in 2022

Zubin Srivastava

```
In [1]: import numpy as np
import pandas as pd
pd.set_option('mode.chained_assignment', None)
import matplotlib.pyplot as plt
from pandasql import sqldf

The dataset used was obtained from Baseball Savant, representing players yearly stats from 2015-2021.
```

```
In [34]: dataset = pd.read_csv('Downloads/stats.csv')
dataset = dataset.rename(columns={"first_name": "first_name"})
dataset = dataset.drop(columns=['Unnamed: 76'])

Out[34]:
```

	last_name	first_name	player_id	year	player_age	b_ab	b_total_pa	b_total_hits	b_single	b_double	...	whiff_percent	swing_perc
0	Colon	Bartolo	112526	2015	42	58	64	8	7	1	...	28.7	...
1	Hunter	Torii	116338	2015	40	521	567	125	81	22	...	23.1	...
2	Ortiz	David	120074	2015	40	528	614	144	70	37	...	23.2	...
3	Rodriguez	Alex	121347	2015	40	523	620	131	75	22	...	32.0	...
4	Ramirez	Aramis	133380	2015	37	475	516	117	68	31	...	17.9	...
...	...	...	...	...	...	...	...	...	...	...	...	...	...
3716	Franco	Wander	677551	2021	20	281	308	81	51	18	...	16.4	...
3717	Duran	Jarren	680776	2021	25	107	112	23	16	3	...	34.7	...
3718	Jeffers	Ryan	680777	2021	24	267	293	53	28	10	...	34.2	...
3719	Miller	Owen	680911	2021	25	191	202	39	27	8	...	26.2	...
3720	Vaughn	Andrew	683734	2021	23	417	469	98	61	22	...	24.6	...

3721 rows x 77 columns

The first goal is to determine the amount of home runs each player will hit per at bat next season. The reason for this is because 2020 was only a 60 game season, I felt that there was a season of less games could affect the accuracy of a model if solely looking for home runs hit in a season. The first step was to create a new column called "hr\_per\_pa" which represented the rate at which players hit a home run per plate appearance.

```
In [3]: dataset['hr_per_pa'] = dataset['b_home_run']/dataset['b_total_pa']

In [4]: dataset = dataset.sort_values(by=['player_id', 'year'], ascending=[True, True])
dataset
```

```
Out[4]:
```

	last_name	first_name	player_id	year	player_age	b_ab	b_total_pa	b_total_hits	b_single	b_double	...	whiff_percent	swing_perc
0	Colon	Bartolo	112526	2015	42	58	64	8	7	1	...	28.7	...
1663	Colon	Bartolo	112526	2016	43	60	65	5	2	2	...	42.7	...
1	Hunter	Torii	116338	2015	40	521	567	125	81	22	...	23.1	...
2	Ortiz	David	120074	2015	40	528	614	144	70	37	...	23.2	...
1664	Ortiz	David	120074	2016	41	537	626	169	82	48	...	19.2	...
...	...	...	...	...	...	...	...	...	...	...	...	...	...
3717	Duran	Jarren	680776	2021	25	107	112	23	16	3	...	34.7	...
3164	Jeffers	Ryan	680777	2020	23	55	62	15	12	0	...	33.3	...
3718	Jeffers	Ryan	680777	2021	24	267	293	53	28	10	...	34.2	...
3719	Miller	Owen	680911	2021	25	191	202	39	27	8	...	26.2	...
3720	Vaughn	Andrew	683734	2021	23	417	469	98	61	22	...	24.6	...

3721 rows x 77 columns

```
In [5]: arr = dataset[['player_id', 'year', 'b_home_run', 'b_total_pa']].to_numpy()
arr
```

```
Out[5]: array([[112526, 2015, 0, 64],
 [112526, 2016, 1, 65],
 [116338, 2015, 22, 567],
 ...,
 [680777, 2021, 14, 293],
 [680911, 2021, 4, 202],
 [683734, 2021, 15, 469]])
```

Next step was to add another column representing next year's hr/pa to train our model. If the season was not 2021 and it was not the player's final season, then we would list their hr rate for the next season, otherwise it would be NULL.

```
In [6]: next_year = []
for row in range(len(arr)):
    if(arr[row, 1] == 2021):
        next_year.append(None)
    elif(arr[row, 0] == arr[row+1, 0]):
        next_year.append((arr[row+1, 2])/arr[row+1, 3])
    else:
        next_year.append(None)
```

The array next year that represents the players hr/pa for their next season is added to the dataframe

```
In [7]: dataset['next_year_hr_per_pa'] = next_year
dataset.iloc[:, [0,1,3,7,-1]]
# Also a good reminder that Bartolo Colon once hit a home run

Out[7]:
```

	last_name	first_name	year	hr_per_pa	next_year_hr_per_pa
0	Colon	Bartolo	2015	0.000000	0.015385
1663	Colon	Bartolo	2016	0.015385	NaN
1	Hunter	Torii	2015	0.038801	NaN
2	Ortiz	David	2015	0.060261	0.060703
1664	Ortiz	David	2016	0.060703	NaN
...	...	...	...	...	...
3717	Duran	Jarren	2021	0.017857	NaN
3164	Jeffers	Ryan	2020	0.048387	0.047782
3718	Jeffers	Ryan	2021	0.047782	NaN
3719	Miller	Owen	2021	0.019802	NaN
3720	Vaughn	Andrew	2021	0.031983	NaN

3721 rows x 5 columns

Next I start preparing to build my model. Model will predict the player's next year hr/pa. I am going to train and test my model on seasons prior to 2021. After it is tested, I will make predictions on data from the 2021 season.

```
In [8]: copy = dataset[(dataset['year'] < 2021) & (dataset['b_total_pa'] >= 100)]
copy = copy.dropna()
copy
```

```
Out[8]:
```

	last_name	first_name	player_id	year	player_age	b_ab	b_total_pa	b_total_hits	b_single	b_double	...	swing_percent	pull_percent
2	Ortiz	David	120074	2015	40	528	614	144	70	37	...	44.7	41
3	Rodriguez	Alex	121347	2015	40	523	620	131	75	22	...	43.9	38
5	Beltre	Adrian	134181	2015	36	567	619	163	109	32	...	48.1	38
1666	Beltre	Adrian	134181	2016	37	583	640	175	111	31	...	48.7	38
2205	Beltre	Adrian	134181	2017	38	340	389	106	66	22	...	46.8	34
...	...	...	...	...	...	...	...	...	...	...	...	...	...
3158	Paredes	Isaac	670623	2020	21	100	108	22	17	4	...	41.9	38
1662	Brosseau	Mike	670712	2019	25	132	142	36	23	7	...	52.6	38
3160	Garcia	Luis	671277	2020	20	134	139	37	29	6	...	49.9	31
3161	Robert	Luis	673357	2020	23	202	227	47	28	8	...	57.6	32
3162	Akiyama	Shogo	673451	2020	32	155	183	38	31	6	...	42.6	28

2191 rows x 78 columns

I originally began with all columns from the data that was downloaded, and I removed columns that were making my mean squared error (MSE) higher, by viewing the feature importance graph below. The columns that are now implemented are the ones that are increasing my accuracy by lowering the MSE. The graph clearly shows that the features xISO and barrel\_batted\_rate are the two most important in increasing accuracy, but there are 13 others that also help the model as well.

```
In [9]: features = copy.iloc[:, [4,9,11,\
                               31,32,\
                               42,44,45,47,50,\
                               61,68,71,75,76]]
features
```

```
Out[9]:
```

	player_age	b_double	b_home_run	xobp	xiso	sweet_spot_percent	barrel_batted_rate	solidcontact_percent	poorlyunder_percent
2	40	37	37	0.385	0.312	34.8	13.1	8.4	24.0
3	40	22	33	0.354	0.246	31.4	10.9	8.1	23.9
5	36	32	18	0.343	0.185	35.7	5.5	7.8	23.5
1666	37	31	32	0.348	0.218	34.2	8.5	5.6	27.9
2205	38	22	17	0.360	0.181	36.1	5.8	6.5	25.5
...	...	...	...	...	...	...	...	...	...
3158	21	4	1	0.268	0.080	35.5	0.0	5.3	28.9
1662	25	7	6	0.265	0.141	31.9	4.3	6.4	34.7
3160	20	6	2	0.293	0.110	27.6	4.8	2.9	10.0
3161	23	8	11	0.293	0.238	36.6	13.0	9.2	24.4
3162	32	6	0	0.355	0.081	34.7	0.8	5.0	19.0

2191 rows x 15 columns

x is my inputs, and y is my outputs, the column "next\_year\_hr\_per\_pa"

```
In [10]: x = np.array(features)
y = np.array(copy.iloc[:, -1])

I created a train_test_split from my x and y data, is 75% training and 25% testing to reduce overfitting
```

```
In [11]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=42)
```

I chose to use a random forest model for predicting each player's next year's hr/pa. I felt a random forest was best because of the multiple decision trees that are used to make the prediction.

```
In [12]: from sklearn.ensemble import RandomForestRegressor
clf = RandomForestRegressor()
clf.fit(x_train, y_train)

Out[12]: RandomForestRegressor()
```

```
In [13]: preds = clf.predict(x_test)
```

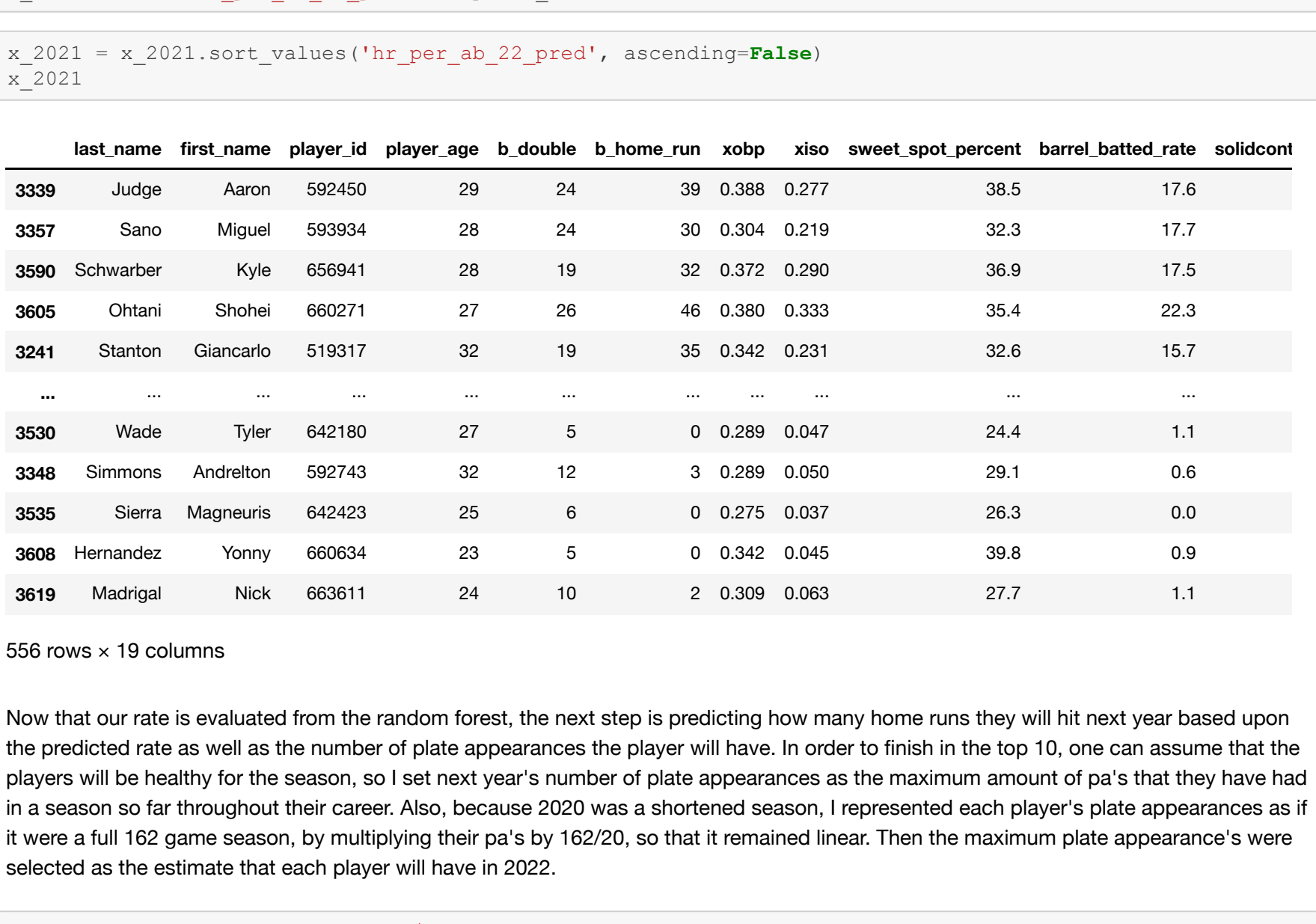
Evaluating feature importance using permutation. Error bars are graphed upon bar graph that represent the standard deviation of each feature for their decrease in accuracy.

```
In [14]: from sklearn.inspection import permutation_importance
feature_names = [feat for feat in features]
result = permutation_importance(clf, x_test, y_test, n_repeats=10, random_state=42, n_jobs=2)
forest_importances = pd.Series(result.importances_mean, index=feature_names)
forest_importances
```

```
Out[14]:
```

player_age	0.001313
b_double	0.005781
b_home_run	0.007743
xobp	0.003490
xiso	0.107544
sweet_spot_percent	0.016038
barrel_batted_rate	0.099880
solidcontact_percent	0.014236
poorlyunder_percent	0.036821
hard_hit_percent	0.037226
pitch_count_breaking	0.011617
swing_percent	0.002037
f_strike_percent	0.005568
popups_percent	0.013486
hr_per_pa	0.011301
dtype:	float64

```
In [15]: plt.figure(figsize=(15,10))
plt.bar(feature_names, forest_importances, yerr = result.importances_std)
plt.xticks(rotation='vertical')
plt.ylabel("Mean decrease in model score")
plt.axhline(y=0, color='r')
plt.title("Feature Importances Using Permutation Model in Evaluating HR Rate")
plt.show()
```



```
In [16]: pd.set_option('display.max_rows', 500)
forest_importances
```

```
Out[16]:
```

player_age	0.001313
b_double	0.005781
b_home_run	0.007743
xobp	0.003490
xiso	0.107544
sweet_spot_percent	0.016038
barrel_batted_rate	0.099880
solidcontact_percent	0.014236
poorlyunder_percent	0.036821
hard_hit_percent	0.037226
pitch_count_breaking	0.011617
swing_percent	0.002037
f_strike_percent	0.005568
popups_percent	0.013486
hr_per_pa	0.011301
dtype:	float64

I used Mean-Squared Error (MSE) to evaluate the quality of the model built upon the features that were used as training data. My MSE is shown below.

```
In [17]: n = x_test.shape[0]
amt = np.sum((preds-y_test)**2)
mse = 1/n*(amt)
mse

Out[17]: 0.00017040566521815564
```

After building the model, it's time run the model on the data from 2021, to predict their hr/pa rate for 2022, and finally predict the top 10 home run hitters for next year

```
In [18]: data_2021 = dataset[(dataset['year'] == 2021)]
data_2021
```

```
Out[18]:
```

	last_name	first_name	player_id	year	player_age	b_ab	b_total_pa	b_total_hits	b_single	b_double	...	swing_percent	pull_percent
3165	Pujols	Albert	405395	2021	41	275	296	65	45	3	...	47.1	41
3166	Cabrera	Miguel	408234	2021	38	472	526	121	90	16	...	49.4	38
3167	Rivera	Rene	425784	2021	38	69	78	16	11	3	...	50.3	41
3168	Wainwright	Adam	425794	2021	40	57	74	7	5	2	...	56.7	38
3169	Molina	Yadier	425877	2021	39	440	473	111	81	19	...	57.1	41
...	...	...	...	...	...	...	...	...	...	...	...	...	...
3716	Franco	Wander	677551	2021	20	281	308	81	51	18	...	49.4	38
3717	Duran	Jarren	680776	2021	25	107	112	23	16	3	...	51.7	27
3718	Jeffers	Ryan	680777	2021	24	267	293	53	28	10	...	46.3	41
3719	Miller	Owen	680911	2021	25	191	202	39	27	8	...	46.9	28
3720	Vaughn	Andrew	683734	2021	23	417	469	98	61	22	...	45.4	38

556 rows x 78 columns

Using the same features that were used to build the random forest

```
In [19]: x_2021 = data_2021.iloc[:, [0,1,2,4,9,11,\
                                   31,32,\
                                   42,44,45,47,50,\
                                   61,68,71,75,76]]

In [20]: x_input = x_2021.iloc[:, 31:].to_numpy()
preds_2022 = clf.predict(x_input)
```

After the model is run on 2022 data, I create a new column on the input data for 2021 called "hr\_per\_ab\_22\_pred" which represents the prediction of hr/pa in 2022

```
In [21]: x_2021.loc[:, 'hr_per_ab_22_pred'] = preds_2022

In [22]: x_2021 = x_2021.sort_values('hr_per_ab_22_pred', ascending=False)
x_2021
```

```
Out[22]:
```

	last_name	first_name	player_id	year	b_double	b_home_run	xobp	xiso	sweet_spot_percent	barrel_batted_rate	solidcont
3339	Sando	Aaron	592450	29	24	39	0.388	0.277	38.3	17.6	
3357	Judge	Miguel	593934	28	24	30	0.304	0.219	32.5	17.7	
3590	Schwarber	Kyle	656941	28	19	32	0.372	0.290	36.9	17.5	
3605	Ohtani	Shohei	660271	27	26	46	0.380	0.333	35.4	22.3	
3241	Stanton	Giancarlo	519317	32	19	35	0.342	0.231	32.6	15.7	
...	...	...	...	...	...	...	...	...	...	...	
3530	Wade	Tyler	642180	27	5	0	0.289	0.047	24.4	1.1	
3348	Simmons	Andreton	592743	32	12	3	0.289	0.050	29.1	0.6	
3535	Sierra	Magneuris	642423	25	6	0	0.275	0.037	26.3	0.0	
3608	Hernandez	Yonny	660634	23	5	0	0.342	0.045	38.0	0.9	
3619	Madrigal	Nick	663611	24	10	2	0.309	0.063	27.7	1.1	

556 rows x 19 columns

Now that our rate is evaluated from the random forest, the next step is predicting how many home runs they will hit next year based upon the predicted rate as well as the number of plate appearances the player will have. In order to finish in the top 10, one can assume that the players will be healthy for the season, so I set next year's number of plate appearances as the maximum amount of pa's that they have had in a season so