

Baseball Value

Matthew Helbig

2019-09-28

Contents

Introduction	1
Getting started	1
Setting things up in SQL	2
“Subsetting” our Batting data	2
Joining Birth Year and Salary	3
Joining Birth Year, Salary, and Batting	3
Getting player value data from Baseball-Reference	5
Split, apply, combine?	8
Reading in our pitching data	9
Subsetting our pitching data	9
Working with our Batting Value Data	13
Creating variables	14
Solving Issues by Subsetting	16
Graphing Average Salary	16

Introduction

The goal of this project is to identify the attributes of the players with the most value and the least value, as well as predicting which players will have the most and least value in the future.

Getting started

The first step we’ll take is to identify which datasets we’ll be working with. The Lahman Database has a lot of good information that we’ll need, namely Player Value data, salary, and dates of birth.

We’ll download that from <http://www.seanlahman.com/baseball-archive/statistics/>

It’s possible to download the database as a .CSV file, and you could technically do everything needed for this project without using SQL, but it’s a good practice to get used to working with databases, especially since a lot of the datasets you’d see in the real world will be bigger than the one we’re currently working with.

Sean Lahman has posted the entire .sql file on his website, so importing it into MySQL is as simple as hitting ‘Data Import’ and loading that database into its own schema.

Once the Lahman database is all settled in its own schema, we want to start to work with that data. There are packages that let you run SQL commands inside R, and there is definitely some value in not having to hop back and forth between R and your SQL environments. However, I’m more comfortable working in MySQL directly, so we won’t worry too much about those R-SQL packages (namely RMySQL).

Setting things up in SQL

The three tables that we’re most interested in right now are “Batting,” “Pitching,” and “Salaries.” The database has a ton of information that’s interesting, but ultimately not super useful for our experiment. We’re going to be mostly concerned with salary information from 2006 to 2016, since a lot of salary data before 2006 is incomplete and we’re going to be running some predictive models on 2017 and 2018 data later.

“Subsetting” our Batting data

We’ll need to create a new table in SQL based on our criteria of only needing 2006 and later (this database ends in 2016).

```
CREATE TABLE Batting_post_2006
LIKE Batting
```

Now we’ve created a table with the same containers as our Batting table, and now we need to load the data into those containers.

```
INSERT INTO Batting_post_2006
SELECT *
FROM Batting;
```

Now we’re going to trim our data from our new table in order to limit ourselves to only 2006 to 2016. This is where SQL is beneficial, since our original dataset of “Batting” is complete safe and untouched.

```
DELETE FROM Batting_post_2006
WHERE yearID < 2006
```

We’re going to do the same thing with our Salaries table, and eventually our Pitching table. We won’t do the Pitching table just yet, but we can follow these same steps in order to set that table up for export to R. For now, we’ll work on the Salaries table:

```
CREATE TABLE Salaries_post_2006
LIKE Salaries
```

```
INSERT INTO Salaries_post_2006
SELECT *
FROM Salaries;
```

```
DELETE FROM Salaries_post_2006
WHERE yearID < 2006
```

Here, we created a new table to put our data in, moved all of our data from the Salaries table into it, and then limited ourselves to data after 2016. Once again, the original Salaries table is untouched.

Joining Birth Year and Salary

One important step that hasn't happened yet is getting birth year statistics from the Master table in the database. This will allow us to create an Age variable in R later, which will be a potentially key attribute for determining who is overvalued and undervalued. In order to add birth year to our Salary data, we'll need to join the two tables.

First, we create a new container for our combined data:

```
CREATE TABLE Player_Salary_Info (  
  AutoNum int NOT NULL AUTO_INCREMENT,  
  playerID varchar(255),  
  birthyear int,  
  yearID int,  
  salary int,  
  PRIMARY KEY (AutoNum)  
);
```

The AutoNum variable is important here, because we need something as the primary key, and if we were to use playerID then we'd have duplicate values between Salary and Master (since the playerIDs are the same between the two tables).

Next, we'll populate our new table with data:

```
INSERT INTO Player_Salary_Info  
SELECT 0, Salaries_post_2006.playerID, Master.birthYear,  
Salaries_post_2006.yearID, Salaries_post_2006.salary  
FROM Salaries_post_2006  
INNER JOIN Master ON Salaries_post_2006.playerID = Master.playerID
```

So now we know that we've managed to get playerID, birthyear, yearID, and salary in the same spot. Our next task will be combining this information with our Batting table, at which point we'll be able to move this new table over to R.

Joining Birth Year, Salary, and Batting

One hiccup we face is that several of our variables for some reason are indicated as strings despite being numeric values. I've encountered issues with this before, where a value of 0 would return an empty string (like " ") instead of actually being 0. To solve this, we can change the string columns into integer columns using the following lines of code:

```
ALTER TABLE example  
ADD newCol int;  
  
UPDATE example SET newCol = CAST(column_name AS UNSIGNED);  
  
ALTER TABLE example  
DROP COLUMN column_name;  
  
ALTER TABLE example  
CHANGE COLUMN newCol column_name int;
```

Now that our data is in a format we can work with that won't give us (hopefully) any errors, we can create and populate our Batting table with salary added. We create the parameters for the new table:

```

CREATE TABLE Batting_Salary (
  AutoNum int NOT NULL AUTO_INCREMENT,
  playerID varchar(255),
  birthyear int,
  yearID int,
  salary int,
  stint int,
  teamID varchar(255),
  lgID varchar(255),
  G int,
  AB int,
  R int,
  H int,
  2B int,
  3B int,
  HR int,
  RBI int,
  SB int,
  CS int,
  BB int,
  SO int,
  GIDP int,
  SF int,
  SH int,
  HBP int,
  IBB int,
  PRIMARY KEY (AutoNum)
);

```

Now we populate our new table:

```

INSERT INTO Batting_Salary
SELECT 0, Batting_post_2006.playerID, Player_Salary_Info.birthyear,
Player_Salary_Info.yearID, Player_Salary_Info.salary, Batting_post_2006.stint,
Batting_post_2006.teamID, Batting_post_2006.lgID, Batting_post_2006.G,
Batting_post_2006.AB, Batting_post_2006.R, Batting_post_2006.H,
Batting_post_2006.2B, Batting_post_2006.3B, Batting_post_2006.HR,
Batting_post_2006.RBI, Batting_post_2006.SB, Batting_post_2006.CS,
Batting_post_2006.BB, Batting_post_2006.SO, Batting_post_2006.GIDP,
Batting_post_2006.SF, Batting_post_2006.SH, Batting_post_2006.HBP,
Batting_post_2006.IBB
FROM Batting_post_2006
INNER JOIN Player_Salary_Info ON Player_Salary_Info.playerID = Batting_post_2006.playerID
AND Player_Salary_Info.yearID = Batting_post_2006.yearID

```

So now we've successfully combined our tables with salary information, birth year information, and batting statistics into one table called Batting_Salary. In the process, we also already shaped the data so that it's only dealing with the years we're concerned with (2006 to 2016). We're just one step away from having our Batter data ready for analysis in R. That step is getting our player value data from Baseball-Reference.

Getting player value data from Baseball-Reference

Our first step is to head to <https://www.baseball-reference.com/data/>, where they keep a daily log of player value data for every player in baseball history. Player value is measured in a lot of different ways, and a constant source of contention in the baseball community is which metrics provide the most accurate measurement of a player's true value.

The player value metric we'll be focusing on the most is Wins Above Replacement (abbreviated as WAR), which measures how many wins a player provided to his team compared to a random player that you could find in the minor leagues. As a rough translation, consider WAR to be the "stock price" of MLB players, as it gives a simple and quick approximation of the value of a player (or stock).

The information we're looking for is all the way at the bottom of the page, under "war_daily_bat.txt." We'll download that, and convert it to .csv using a spreadsheet program (I'm using Open Office). Just follow the steps for converting everything, and we'll save it to our active folder.

Now we need to load this data into SQL. However, if we look at this data in R:

```
Daily_bat <- read.csv("war_daily_bat.csv")
names(Daily_bat)
```

```
## [1] "name_common"      "age"              "mlb_ID"
## [4] "player_ID"        "year_ID"          "team_ID"
## [7] "stint_ID"         "lg_ID"            "PA"
## [10] "G"                "Inn"              "runs_bat"
## [13] "runs_br"          "runs_dp"          "runs_field"
## [16] "runs_infield"     "runs_outfield"    "runs_catcher"
## [19] "runs_good_plays"  "runs_defense"     "runs_position"
## [22] "runs_position_p"  "runs_replacement" "runs_above_rep"
## [25] "runs_above_avg"   "runs_above_avg_off" "runs_above_avg_def"
## [28] "WAA"              "WAA_off"          "WAA_def"
## [31] "WAR"              "WAR_def"          "WAR_off"
## [34] "WAR_rep"          "salary"           "pitcher"
## [37] "teamRpG"          "oppRpG"           "oppRpPA_rep"
## [40] "oppRpG_rep"       "pyth_exponent"    "pyth_exponent_rep"
## [43] "waa_win_perc"     "waa_win_perc_off" "waa_win_perc_def"
## [46] "waa_win_perc_rep" "OPS_plus"         "TOB_lg"
## [49] "TB_lg"
```

We can see that we have a ton of variables that we don't necessarily want. Remember when we created a container in SQL above? If we were to try to read this into SQL now, we'd have to create a container for every one of these variables. We don't necessarily need "oppRpG," which measures how many runs per game that player's opponent scored. Similarly, we don't need a lot of other variables. We're going to limit ourselves to just the ones we want to keep by running the following (Note: you'll need the "dplyr" library for this function):

```
kept_Columns = select(Daily_bat, 1, 2, 4, 5, 6, 31, 32, 33, 36, 47)
```

Which is going to keep only the columns we're interested in, namely:

```
names(kept_Columns)
```

```
## [1] "name_common" "age"          "player_ID"    "year_ID"      "team_ID"
## [6] "WAR"         "WAR_def"      "WAR_off"      "pitcher"      "OPS_plus"
```

Next we're going to subset our data to keep it within the year range we're investigating (2006 to 2016), as well as limiting our data to be only batters, since we really aren't interested in how pitchers hit. Pitchers aren't paid based on their ability to hit, so batting statistics for pitchers aren't relevant to what we're studying.

```
years_pitchers_Subset <- subset(kept_Columns, 2005 < year_ID & year_ID < 2017 & pitcher == "N")
```

Finally, we'll remove the pitcher column, since all of our kept players are going to be batters:

```
Batting_value <- select(years_pitchers_Subset, -9)
```

After all this, we're ready to export this data from R and into SQL so we can join it with our Batting_Salary table.

```
write.csv(Batting_value, "Batting_Value.csv")
```

We create our container in SQL by running the following:

```
CREATE TABLE Batting_Value (  
  AutoNum int NOT NULL AUTO_INCREMENT,  
  name_common varchar(255),  
  age int,  
  player_ID varchar(255),  
  year_ID int,  
  team_ID varchar(255),  
  WAR DECIMAL(4,2),  
  WAR_def DECIMAL(4,2),  
  WAR_off DECIMAL(4,2),  
  OPS_plus int,  
  PRIMARY KEY (AutoNum)  
);
```

Then we'll load our data into the container using the following command:

```
LOAD DATA LOCAL INFILE '/filepath/Batting_Value.csv'  
INTO TABLE Batting_Value FIELDS TERMINATED BY ','  
ENCLOSED BY '"' LINES TERMINATED BY '\n'  
IGNORE 1 LINES
```

Now we want to combine this data with our Batting_Salary data so we can have our final dataset to work with in R.

We'll create the container we'll use for our combined data:

```
CREATE TABLE Value_Salary (  
  AutoNum int NOT NULL AUTO_INCREMENT,  
  playerID varchar(255),  
  name_common varchar(255),  
  yearID int,  
  age int,  
  salary int,  
  stint int,  
  teamID varchar(255),
```

```

lgID varchar(255),
G int,
AB int,
R int,
H int,
2B int,
3B int,
HR int,
RBI int,
SB int,
CS int,
BB int,
SO int,
GIDP int,
SF int,
SH int,
HBP int,
IBB int,
WAR DECIMAL(4,2),
WAR_def DECIMAL(4,2),
WAR_Off DECIMAL(4,2),
OPS_plus int,
PRIMARY KEY (AutoNum)
);

```

Then we load in our data:

```

INSERT INTO Value_Salary
SELECT DISTINCT 0,
Batting_Salary.playerID,
Batting_Value.name_common,
Batting_Salary.yearID,
Batting_Value.age,
Batting_Salary.salary,
Batting_Salary.stint,
Batting_Salary.teamID,
Batting_Salary.lgID,
Batting_Salary.G,
Batting_Salary.AB,
Batting_Salary.R,
Batting_Salary.H,
Batting_Salary.2B,
Batting_Salary.3B,
Batting_Salary.HR,
Batting_Salary.RBI,
Batting_Salary.SB,
Batting_Salary.CS,
Batting_Salary.BB,
Batting_Salary.SO,
Batting_Salary.GIDP,
Batting_Salary.SF,
Batting_Salary.SH,
Batting_Salary.HBP,

```

```

Batting_Salary.IBB,
Batting_Value.WAR,
Batting_Value.WAR_def,
Batting_Value.WAR_off,
Batting_Value.OPS_plus
FROM Batting_Salary
INNER JOIN Batting_Value ON Batting_Salary.playerID = Batting_Value.player_ID AND Batting_Salary.yearID

```

As a note, this table works very well until it comes to the players who were traded in the middle of the season. When a player hasn't been traded, there is one row per year. However, when a player has been traded, we would expect to see two columns that year (one for the old team, one for the new team). However, we see four:

##	AutoNum	playerID	name_common	yearID	age	salary	stint
## 1	1	abadan01	Andy Abad	2006	33	327000	1
## 2	2	abercro01	Reggie Abercrombie	2006	25	327000	1
## 3	3	abreubo01	Bobby Abreu	2006	32	13600000	1
## 4	4	abreubo01	Bobby Abreu	2006	32	13600000	1
## 5	5	abreubo01	Bobby Abreu	2006	32	13600000	2
## 6	6	abreubo01	Bobby Abreu	2006	32	13600000	2

Split, apply, combine?

We need to remove those duplicates, and one very common way of doing so is via a process often used in data analysis called “Split, apply, combine.”

This process involves splitting the data into subsets, applying a function to each subset, and then combining the data back together. Which is what we want to do here. We want to remove the duplicate players from every year, but we don't want to remove duplicates from the dataset as a whole. For example, there are multiple “Bobby Abreu” rows in our 2006 data, which we would like to reduce to a single row. However, we don't want to remove all of the “Bobby Abreu” rows from our master dataset, because that would remove Bobby Abreu's stats for 2007, 2008, and so on.

We could subset our data by year, write a function that removes duplicates for each year, and then join all of our subsets back together. Or, we could use a library that we've already called, the ‘dplyr’ library, and remove the duplicates by year from our main dataset.

If we execute the following command:

```

Batting_Value_Salary <- Pre_Batting_Value_Salary %>%
  distinct(playerID, yearID, .keep_all = TRUE)

```

We can see that we've successfully removed the duplicates:

##	AutoNum	playerID	name_common	yearID	age	salary	stint
## 1	1	abadan01	Andy Abad	2006	33	327000	1
## 2	2	abercro01	Reggie Abercrombie	2006	25	327000	1
## 3	3	abreubo01	Bobby Abreu	2006	32	13600000	1
## 4	7	adamsru01	Russ Adams	2006	25	343000	1
## 5	8	aguilch01	Chris Aguila	2006	27	327000	1
## 6	9	alfoned01	Edgardo Alfonzo	2006	32	8000000	1

This is a nice way of using pre-built libraries to limit the amount of code that you need to write. The dplyr library has a ton of functions like this, and we'll make use of them when we further analyze our data in R.

Reading in our pitching data

Our next step will be reading in our pitching data.

Subsetting our pitching data

First we create our container in SQL:

```
CREATE TABLE Pitching_post_2006
LIKE Pitching
```

Then we load our data into our table:

```
INSERT INTO Pitching_post_2006
SELECT *
FROM Pitching;
```

Then we do what we did with Batting and Salary and limit our data to just 2006-2016:

```
DELETE FROM Pitching_post_2006
WHERE yearID < 2006
```

We create our Pitching_Salary table:

```
CREATE TABLE Pitching_Salary (
    AutoNum int NOT NULL AUTO_INCREMENT,
    playerID varchar(255),
    birthyear int,
    yearID int,
    salary int,
    stint int,
    teamID varchar(255),
    lgID varchar(255),
    W int,
    L int,
    G int,
    GS int,
    CG int,
    SHO int,
    SV int,
    IPouts int,
    H int,
    ER int,
    HR int,
    BB int,
    SO int,
    ERA float,
    BK int,
    R int,
    PRIMARY KEY (AutoNum)
);
```

Then we load our Pitching and Salary data into our new table:

```
INSERT INTO Pitching_Salary
SELECT 0,
Pitching_post_2006.playerID,
Player_Salary_Info.birthyear,
Player_Salary_Info.yearID,
Player_Salary_Info.salary,
Pitching_post_2006.stint,
Pitching_post_2006.teamID,
Pitching_post_2006.lgID,
Pitching_post_2006.W,
Pitching_post_2006.L,
Pitching_post_2006.G,
Pitching_post_2006.GS,
Pitching_post_2006.CG,
Pitching_post_2006.SH0,
Pitching_post_2006.SV,
Pitching_post_2006.IPouts,
Pitching_post_2006.H,
Pitching_post_2006.ER,
Pitching_post_2006.HR,
Pitching_post_2006.BB,
Pitching_post_2006.SO,
Pitching_post_2006.ERA,
Pitching_post_2006.BK,
Pitching_post_2006.R
FROM Pitching_post_2006
INNER JOIN Player_Salary_Info ON Player_Salary_Info.playerID = Pitching_post_2006.playerID
AND Player_Salary_Info.yearID = Pitching_post_2006.yearID
```

So now we're at the point where we need to go back to Baseball-Reference and get our value data. For batters, it was called "war_daily_bat," and for pitchers it's "war_daily_pitch." We take a look at the data in R, and see the following:

```
Daily_pitch <- read.csv("war_daily_pitch.csv")
names(Daily_pitch)

## [1] "name_common"      "age"
## [3] "mlb_ID"           "player_ID"
## [5] "year_ID"          "team_ID"
## [7] "stint_ID"         "lg_ID"
## [9] "G"                "GS"
## [11] "IPouts"           "IPouts_start"
## [13] "IPouts_relief"    "RA"
## [15] "xRA"              "xRA_sprp_adj"
## [17] "xRA_def_pitcher" "PPF"
## [19] "PPF_custom"      "xRA_final"
## [21] "BIP"             "BIP_perc"
## [23] "RS_def_total"    "runs_above_avg"
## [25] "runs_above_avg_adj" "runs_above_rep"
## [27] "RpO_replacement" "GR_leverage_index_avg"
## [29] "WAR"             "salary"
```

```
## [31] "teamRpG"          "oppRpG"
## [33] "pyth_exponent"    "waa_win_perc"
## [35] "WAA"              "WAA_adj"
## [37] "oppRpG_rep"       "pyth_exponent_rep"
## [39] "waa_win_perc_rep" "WAR_rep"
## [41] "ERA_plus"         "ER_lg"
```

Limit ourselves to the names we want, and show the new column names

```
pitch_kept_columns = select(Daily_pitch, 1, 2, 4, 5, 6, 29, 41)
names(pitch_kept_columns)
```

```
## [1] "name_common" "age"          "player_ID"    "year_ID"      "team_ID"
## [6] "WAR"          "ERA_plus"
```

Subset to limit ourselves to the years we want, and write it to csv so we can export it to SQL to combine with our Pitching Salary table

```
Pitching_value <- subset(pitch_kept_columns, 2005 < year_ID & year_ID < 2017)
write.csv(Pitching_value, "Pitching_Value.csv")
```

Create our pitching value container in SQL:

```
CREATE TABLE Pitching_Value (
  AutoNum int NOT NULL AUTO_INCREMENT,
  name_common varchar(255),
  age int,
  player_ID varchar(255),
  year_ID int,
  team_ID varchar(255),
  WAR DECIMAL(4,2),
  ERA_plus int,
  PRIMARY KEY (AutoNum)
);
```

Then load pitching value with data:

```
LOAD DATA LOCAL INFILE '/filepath/Pitching_Value.csv'
INTO TABLE Pitching_Value FIELDS TERMINATED BY ','
ENCLOSED BY '"' LINES TERMINATED BY '\n'
IGNORE 1 LINES
```

Now we want to combine Pitching_Value with Pitching_Salary. We start by creating the container for our data:

```
CREATE TABLE Pitching_Value_Salary (
  AutoNum int NOT NULL AUTO_INCREMENT,
  playerID varchar(255),
  name_common varchar(255),
  yearID int,
```

```

    age int,
    salary int,
    stint int,
    teamID varchar(255),
    lgID varchar(255),
    W int,
    L int,
    G int,
    GS int,
    CG int,
    SHO int,
    SV int,
    IPouts int,
    H int,
    ER int,
    HR int,
    BB int,
    SO int,
    ERA float,
    BK int,
    R int,
    WAR DECIMAL(4,2),
    ERA_plus int,
    PRIMARY KEY (AutoNum)
);

```

Then we load in our data:

```

INSERT INTO Pitching_Value_Salary
SELECT DISTINCT 0,
Pitching_Salary.playerID,
Pitching_Value.name_common,
Pitching_Salary.yearID,
Pitching_Value.age,
Pitching_Salary.salary,
Pitching_Salary.stint,
Pitching_Salary.teamID,
Pitching_Salary.lgID,
Pitching_Salary.W,
Pitching_Salary.L,
Pitching_Salary.G,
Pitching_Salary.GS,
Pitching_Salary.CG,
Pitching_Salary.SHO,
Pitching_Salary.SV,
Pitching_Salary.IPouts,
Pitching_Salary.H,
Pitching_Salary.ER,
Pitching_Salary.HR,
Pitching_Salary.BB,
Pitching_Salary.SO,
Pitching_Salary.ERA,
Pitching_Salary.BK,

```

```
Pitching_Salary.R,
Pitching_Value.WAR,
Pitching_Value.ERA_plus
FROM Pitching_Salary
INNER JOIN Pitching_Value ON Pitching_Salary.playerID = Pitching_Value.player_ID AND Pitching_Salary.yea
```

We have the same problem with duplicates that we did with our Batting data. Taking a look at the following, we can see those duplicates:

```
Pre_Pitching_Value_Salary <- read.csv("Pitching_Value_Salary.csv")
show_duplicates2 <- subset(Pre_Pitching_Value_Salary, name_common == "Jeremy Accardo")
head(show_duplicates2[,1:7])
```

##	AutoNum	playerID	name_common	yearID	age	salary	stint
## 8	8	accarje01	Jeremy Accardo	2006	24	330000	1
## 9	9	accarje01	Jeremy Accardo	2006	24	330000	2
## 10	10	accarje01	Jeremy Accardo	2006	24	330000	1
## 11	11	accarje01	Jeremy Accardo	2006	24	330000	2
## 12	12	accarje01	Jeremy Accardo	2007	25	392200	1
## 13	13	accarje01	Jeremy Accardo	2008	26	392200	1

However, the same distinct function we used from Batting will work here as well:

```
Pitching_Value_Salary <- Pre_Pitching_Value_Salary %>%
  distinct(playerID, yearID, .keep_all = TRUE)
```

Taking a look to see if those duplicates are gone:

```
show_duplicates3 <- subset(Pitching_Value_Salary, name_common == "Jeremy Accardo")
head(show_duplicates3[,1:7])
```

##	AutoNum	playerID	name_common	yearID	age	salary	stint
## 8	8	accarje01	Jeremy Accardo	2006	24	330000	1
## 9	12	accarje01	Jeremy Accardo	2007	25	392200	1
## 10	13	accarje01	Jeremy Accardo	2008	26	392200	1
## 11	14	accarje01	Jeremy Accardo	2010	28	1080000	1
## 12	15	accarje01	Jeremy Accardo	2011	29	1080000	1

There we go! Our Batting Value and Pitching Value data is all ready for analysis in R. It took a lot of behind the scenes work, but it left us with some very high quality data on which to test our hypotheses.

Working with our Batting Value Data

Now that we have everything loaded in and set up properly, we're going to begin the actual analysis step of this project. We spent a ton of time loading in the necessary variable in order to work with our data, but there are still more new variables that will help us draw even deeper conclusions.

Creating variables

Cost-per-win

Since we have wins above replacement data and salary data, we can create a new variable called “cost-per-win” that will let us know how expensive a win is on the open market. If a team were to go out and sign a player, the “cost-per-win” would be how much it would cost, on average, to buy 1 win above replacement. That data has been collected online, and we’re going to use the data from “<https://batflipsandnerds.com/2018/06/20/the-modern-myths-of-baseball-the-cost-of-a-win/>”

These values change every year, so we’ll need to make sure that the specific cost-per-win is applied to the appropriate year. In order to do this, we’ll need to do another split-apply-combine on our data.

First we’ll split the data by year:

```
year <- split(Batting_Value_Salary, Batting_Value_Salary$yearID)
```

Then we’ll apply the appropriate cost-per-win to each year:

```
year$`2006`$cost_per_win <- 4600000
year$`2007`$cost_per_win <- 5300000
year$`2008`$cost_per_win <- 5600000
year$`2009`$cost_per_win <- 5700000
year$`2010`$cost_per_win <- 5800000
year$`2011`$cost_per_win <- 6400000
year$`2012`$cost_per_win <- 6900000
year$`2013`$cost_per_win <- 7500000
year$`2014`$cost_per_win <- 7700000
year$`2015`$cost_per_win <- 8700000
year$`2016`$cost_per_win <- 9600000
```

Finally, we’ll combine our data back together now that we’ve split it by year and applied the appropriate cost-per-win:

```
Batting_Value_Salary <- rbind(year$`2006`, year$`2007`, year$`2008`, year$`2009`, year$`2010`, year$`2011`, year$`2012`, year$`2013`, year$`2014`, year$`2015`, year$`2016`)
```

Player Money Value

The next variable we’re going to create is going to be “Player Money Value.” This is going to capture how much value, in dollars, a player provided that season. We’ll do this by multiplying a player’s wins above replacement by the cost of a win above replacement. This number will essentially be how much a player was worth if they had been available on the open market and received fair market value. This is how we’ll do it:

```
Batting_Value_Salary$player_money_value <- (Batting_Value_Salary$WAR * Batting_Value_Salary$cost_per_win)
```

Value Plus

It’s important to de-contextualize a player’s value. For example, if we said that Joe Random provided \$6 million in value, we wouldn’t know off the top of our heads if that was good or bad. We’d need to know how much value everyone else was providing. “Value Plus” is a statistic that allows us to do that. It normalizes the value a player provides, where 100 is average. If a player has a Value Plus of 120, it means that they provided 20% more value than average. We’ll calculate this by doing the following:

```
Batting_Value_Salary$value_plus <- ((Batting_Value_Salary$player_money_value / mean(Batting_Value_Salary$player_money_value)) - Batting_Value_Salary$player_money_value)
```

Excess value

The next variable we'll need is "Excess value." This will measure the difference between how much a player was worth and how much they actually made. It's worth noting that for some players, especially good players with low salaries (i.e. players still on their rookie contracts), this number will be quite high. For some players, this number could be in the negatives, indicating that they cost their team money with their play. We'll calculate this as follows:

```
Batting_Value_Salary$excess_value <- (Batting_Value_Salary$player_money_value - Batting_Value_Salary$salary)
```

Excess value plus

The final variable we'll be working with right now is "Excess Value Plus." This will be a normalized version of "Excess value" in order to give a more accurate and context-free representation of the excess value for each player. While knowing how valuable a player was with respect to other players, it's even more valuable to know how much excess value they provided relative to other players. This is very similar to our "Value Plus" statistic, except this new variable also incorporates a player's salary into it. Here's how we'll calculate it:

```
Batting_Value_Salary$excess_value_plus <- ((Batting_Value_Salary$excess_value / mean(Batting_Value_Salary$excess_value)) - Batting_Value_Salary$excess_value)
```

Batting average, On-Base Percentage, Slugging Percentage, OPS, and ISO

We've mostly created sabermetric statistics thus far, but in order to identify potential undervalued and overvalued players, it will help to use some traditional statistics as well.

Batting average measures how many hits a player gets per at-bat, on-base percentage measures how often a player gets on-base per at-bat, and slugging percentage measures how many bases a player obtains per at-bat (essentially, how hard does he hit the ball).

We'll calculate batting average as follows:

```
pre_batting_average <- (Batting_Value_Salary$H)/(Batting_Value_Salary$AB)
Batting_Value_Salary$batting_average <- round(pre_batting_average, digits = 3)
```

Next, we'll calculate on-base percentage:

```
pre_on_base_percentage <- (Batting_Value_Salary$H + Batting_Value_Salary$BB + Batting_Value_Salary$HBP)/(Batting_Value_Salary$AB)
Batting_Value_Salary$on_base_percentage <- round(pre_on_base_percentage, digits = 3)
```

Then we'll calculate slugging percentage:

```
pre_slugging_percentage <- ((Batting_Value_Salary$H + Batting_Value_Salary$2B + Batting_Value_Salary$3B)/(Batting_Value_Salary$AB))
Batting_Value_Salary$slugging_percentage <- round(pre_slugging_percentage, digits = 3)
```

OPS, or on-base plus slugging, is a statistics that, simply put, combines on-base-percentage and slugging percentage. It's a quick way to see how often a guy got on baseball, and how hard he hit the ball to get on base.

```
Batting_Value_Salary$OPS <- (Batting_Value_Salary$on_base_percentage) + (Batting_Value_Salary$slugging_
```

ISO, or isolated power, is another quick metric to determine how powerful someone was at the plate. The higher the ISO, the more bases someone accumulated relative to their batting average. It's not super helpful as a predictive statistic, but it will give us another tool to work with when we're doing our player similarity evaluations.

```
Batting_Value_Salary$ISO <- (Batting_Value_Salary$slugging_percentage) - (Batting_Value_Salary$batting_
```

Strikeout percentage and walk percentage

The final two statistics we'll create for the time being are going to measure plate discipline. Generally speaking, a low strikeout percentage and a high walk percentage are positive predictors of future success. In order to identify potential breakout players, we'll want to quantify these statistics to use for future evaluation.

```
pre_K_percentage <- ((Batting_Value_Salary$SO)/(Batting_Value_Salary$AB + Batting_Value_Salary$BB + Bat
Batting_Value_Salary$K_percentage <- round(pre_K_percentage, digits = 2)
```

```
pre_BB_percentage <- ((Batting_Value_Salary$BB)/(Batting_Value_Salary$AB + Batting_Value_Salary$BB + Ba
Batting_Value_Salary$BB_percentage <- round(pre_BB_percentage, digits = 2)
```

Solving Issues by Subsetting

Now that we have all of our variables created, we seem like we're just about set to start getting to the meat of the project, and determining which players are the most undervalued and which are the most overvalued.

However, some quick graphs will highlight an issue that we face, and this issue isn't so much with our data so much as it is with an underlying fact about Major League Baseball.

Graphing Average Salary

Since our main focus will be on the amount of excess value a player provided, and excess value is comprised of player value and salary, it will be key to focus on the average salaries for players in different age groups. This will demonstrate a bit where our problem lies.

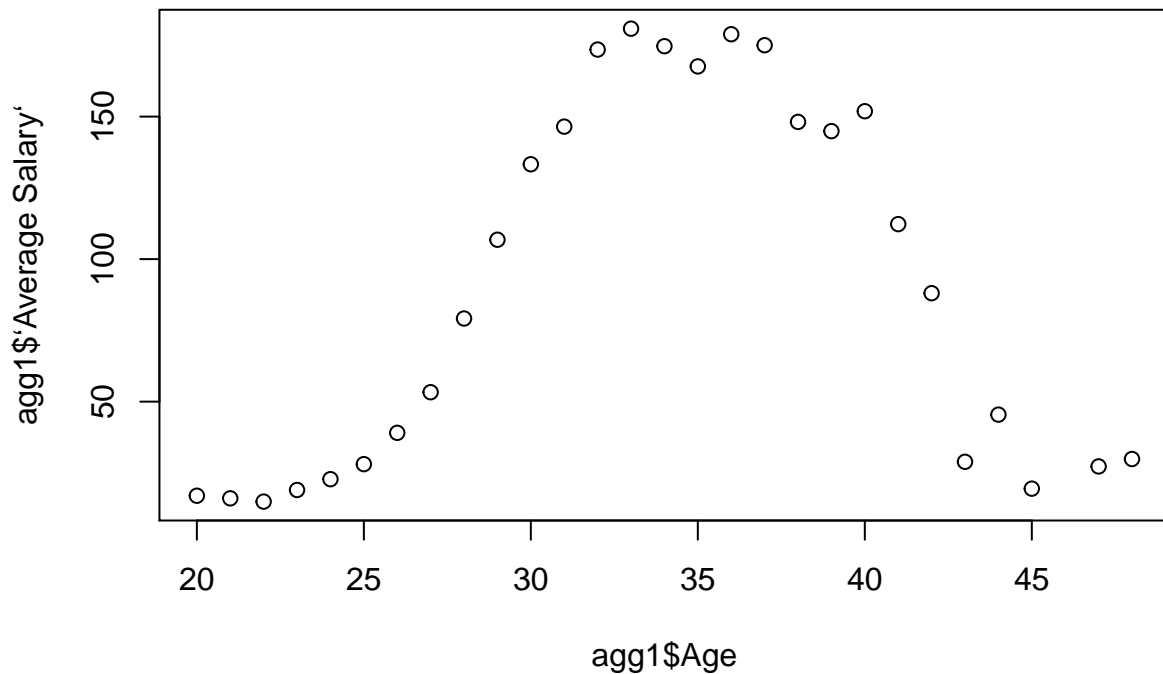
First, we'll create a Salary+ metric that will normalize salary data relative to the average salary.

```
Batting_Value_Salary$salary_plus <- ((Batting_Value_Salary$salary / mean(Batting_Value_Salary$salary)*1
```

Next, we're going to calculate the mean salary by age, and then plot that to get a graphical representation of how salary changes as a player gets older.

```
agg1 <- aggregate(x = Batting_Value_Salary$salary_plus, by = list(Batting_Value_Salary$age), FUN = mean,
names(agg1)[1] <- "Age"
names(agg1)[2] <- "Average Salary"

plot(agg1$Age, agg1$`Average Salary`)
```

Notes to self

Tomorrow, 4/16, I want to move a lot of the work from the test2 table into the actual table we'll be working with. I also want to create a table of Excess Value Plus and Salary Plus, and quantify the results we're seeing from our line graph. Once we do that, we can justify our age subset, and then start getting to work creating statistics like batting average, on base percentage, slugging percentage, and OPS. Once we have all of our statistics created (which may be before or after we subset), we can do a big multiple regression and see which statistics impact Excess Value Plus the most. I'm not sure what to do with "player_money_value" as a column, because it could be valuable. I'll keep that in the back of my mind.