

MSD 2019 Final Project

A replication and extension of Wage disparity and team productivity: evidence from Major League Baseball by Craig A. Depken II, 1999

Forrest Hofmann (fhh2112) & Sagar Lal (sl3946)

2019-05-07 17:12:45

Contents

Introduction	1
Problem Description	1
Motivation	1
Data Source	1
Reproduction	1
Reproduction Code	1
Reproduction Notes	3
Reproduction Analysis	3
Extension	3
Extension Code	3
Extension Notes	12
Extension Analysis	12
Postface	12

Introduction

Problem Description

Motivation

Data Source

Reproduction

Reproduction Code

```
teams <- read_csv(here('data/teams.csv'))
salaries <- read_csv(here('data/salaries.csv'))

teams$WSWin <- as.logical(teams$WSWin == 'Y')
teams <- teams %>%
  filter(1985 <= yearID & yearID <= 2016) %>%
  mutate(winPercentage = W / (W + L) * 1000) %>%
  filter(yearID != 1987) %>%
  filter(teamID != 'TEX')
```

```

salaries <- salaries %>%
  filter(1985 <= yearID & yearID <= 2016) %>%
  mutate(salaryMil = salary / 1000000) %>%
  filter(yearID != 1987) %>%
  filter(teamID != 'TEX')

#Clearly missing lots of data for this datapoint hence its removal https://www.baseball-reference.com/t

teams <- teams %>%
  inner_join(salaries) %>%
  group_by(yearID, teamID, G, W, L, WSWin, winPercentage) %>%
  summarize(totalSalaryMil = sum(salaryMil))

salaries <- salaries %>%
  inner_join(teams) %>%
  mutate(salaryShare = salaryMil / totalSalaryMil * 100) %>%
  mutate(salaryShareSquared = salaryShare ^ 2) %>%
  select(yearID, teamID, playerID, salary, salaryShare, salaryShareSquared)

teams <- teams %>%
  inner_join(salaries) %>%
  group_by(yearID, teamID, G, W, L, winPercentage, WSWin, totalSalaryMil) %>%
  summarize(HHI = sum(salaryShareSquared))

teams_old <- teams %>%
  filter(1985 <= yearID & yearID <= 1998) %>%
  mutate(normalizedYear = yearID - 1985)

salaries_old <- salaries %>%
  filter(1985 <= yearID & yearID <= 1998)

summary(teams_old$winPercentage)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   327.2   456.8   496.9   500.1   543.2   703.7
sd(teams_old$winPercentage)

## [1] 67.34053
summary(teams_old$totalSalaryMil)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    4.613   14.217   23.655   26.220   37.022   72.356
sd(teams_old$totalSalaryMil)

## [1] 14.00647
summary(teams_old$HHI)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    427.5   666.8   754.7   801.5   876.3   2158.3
sd(teams_old$HHI)

## [1] 220.2525

```

```
linear_fixed_old <- lm(formula = winPercentage ~ totalSalaryMil + HHI +
                      normalizedYear + teamID + 0,
                      data = teams_old)
summary(linear_fixed_old)$coefficients[1:3,]

##              Estimate Std. Error   t value    Pr(>|t|)
## totalSalaryMil  1.96115149 0.48091620  4.077948 5.803436e-05
## HHI             -0.04611478 0.02028827 -2.272977 2.372293e-02
## normalizedYear -4.08667575 1.76297250 -2.318060 2.110738e-02

linear_random_old <- lm(formula = winPercentage ~ totalSalaryMil + HHI + normalizedYear,
                      data = teams_old)
summary(linear_random_old)$coefficients[1:4,]

##              Estimate Std. Error   t value    Pr(>|t|)
## (Intercept)    515.9993052 15.18582553 33.979009 1.416422e-110
## totalSalaryMil   2.1067707  0.41153306  5.119323 5.180395e-07
## HHI             -0.0469241  0.01853488 -2.531665 1.180894e-02
## normalizedYear  -4.7782865  1.57044613 -3.042630 2.530509e-03
```

Reproduction Notes

- original author did not describe how time fixed effects are accounted for (across expansion periods or every year)
- no discussion about limiting to 25 man roster vs 40 man roster
- no discussion of cut players, traded players
- no discussion of signing bonuses

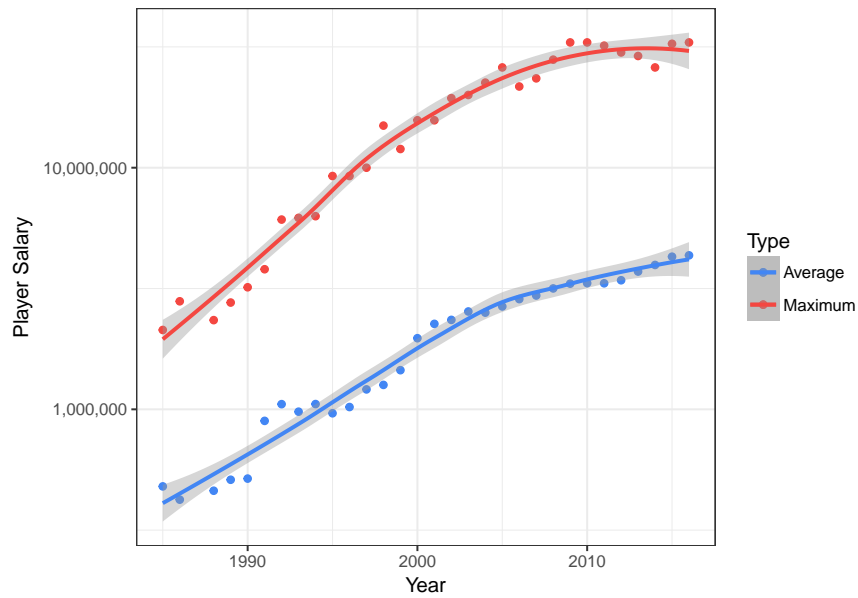
Reproduction Analysis

Extension

Extension Code

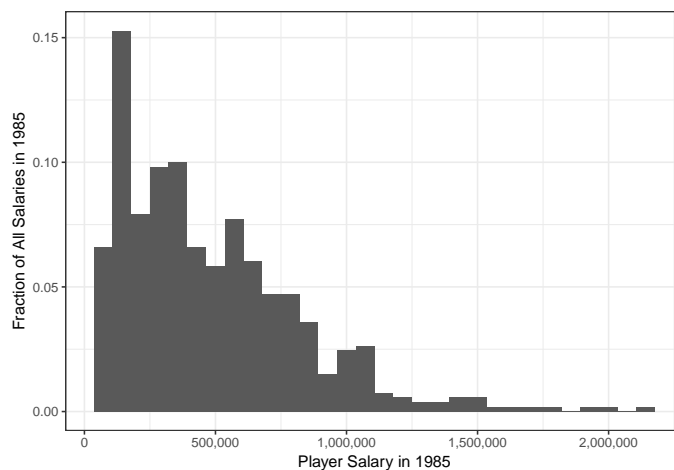
```
salary_vs_time <- salaries %>%
  group_by(yearID) %>%
  summarize(avg = mean(salary), max = max(salary))

ggplot(data = salary_vs_time) +
  geom_point(aes(x = yearID, y = avg, color = 'Average')) +
  geom_smooth(aes(x = yearID, y = avg, color = 'Average')) +
  geom_point(aes(x = yearID, y = max, color = 'Maximum')) +
  geom_smooth(aes(x = yearID, y = max, color = 'Maximum')) +
  scale_color_manual(values = c('#4286f4', '#f44741')) +
  scale_y_log10(labels = comma) +
  labs(color = 'Type') +
  xlab('Year') +
  ylab('Player Salary')
```

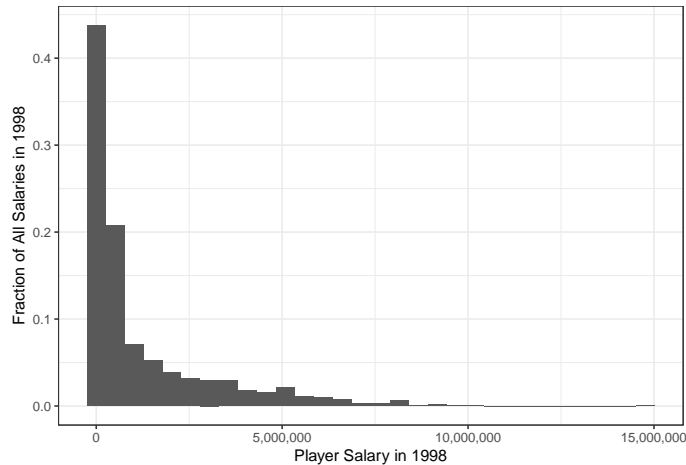


```
salaries_1985 <- filter(salaries, yearID == 1985)
salaries_1998 <- filter(salaries, yearID == 1998)
salaries_2016 <- filter(salaries, yearID == 2016)

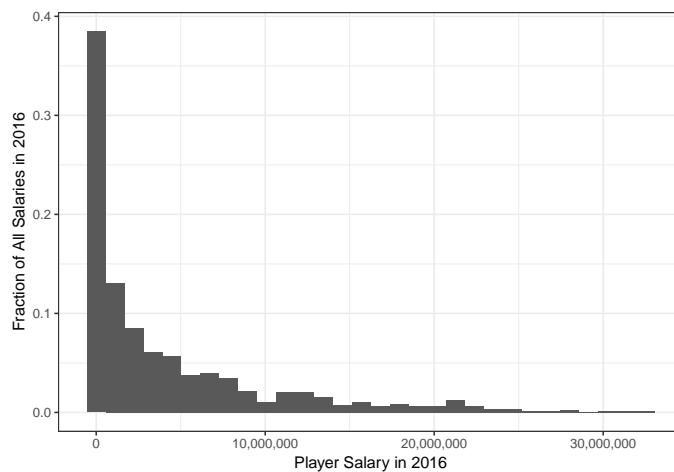
ggplot(data = salaries_1985) +
  geom_histogram(aes(x = salary, y = (..count..) / sum(..count..))) +
  scale_x_continuous(labels = comma) +
  xlab('Player Salary in 1985') +
  ylab('Fraction of All Salaries in 1985')
```



```
ggplot(data = salaries_1998) +
  geom_histogram(aes(x = salary, y = (..count..) / sum(..count..))) +
  scale_x_continuous(labels = comma) +
  xlab('Player Salary in 1998') +
  ylab('Fraction of All Salaries in 1998')
```



```
ggplot(data = salaries_2016) +
  geom_histogram(aes(x = salary, y = (..count..) / sum(..count..))) +
  scale_x_continuous(labels = comma) +
  xlab('Player Salary in 2016') +
  ylab('Fraction of All Salaries in 2016')
```

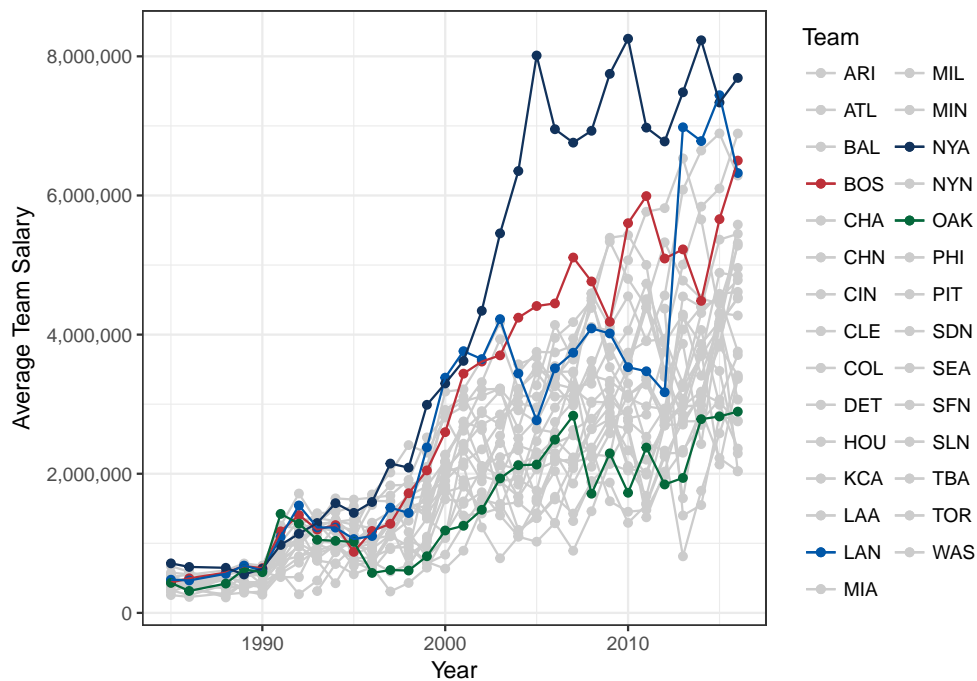


```
current_teamIDs <- c('ARI', 'ATL', 'BAL', 'BOS', 'CHA', 'CHN', 'CIN', 'CLE', 'COL', 'DET',
                    'HOU', 'KCA', 'LAA', 'LAN', 'MIA', 'MIL', 'MIN', 'NYA', 'NYN', 'OAK',
                    'PHI', 'PIT', 'SDN', 'SEA', 'SFN', 'SLN', 'TBA', 'TEX', 'TOR', 'WAS')
team_colors <- c('#cccccc', '#cccccc', '#cccccc', '#BD3039', '#cccccc',
                '#cccccc', '#cccccc', '#cccccc', '#cccccc', '#cccccc',
                '#cccccc', '#cccccc', '#cccccc', '#0157a8', '#cccccc',
                '#cccccc', '#cccccc', '#11325b', '#cccccc', '#04683b',
                '#cccccc', '#cccccc', '#cccccc', '#cccccc', '#cccccc',
                '#cccccc', '#cccccc', '#cccccc', '#cccccc', '#cccccc')
colored_teamIDs <- c('BOS', 'LAN', 'NYA', 'OAK')

team_salary_vs_time <- salaries %>%
  filter(teamID %in% current_teamIDs) %>%
  group_by(yearID, teamID) %>%
  summarize(avg = mean(salary)) %>%
  mutate(flag = teamID %in% colored_teamIDs)
```

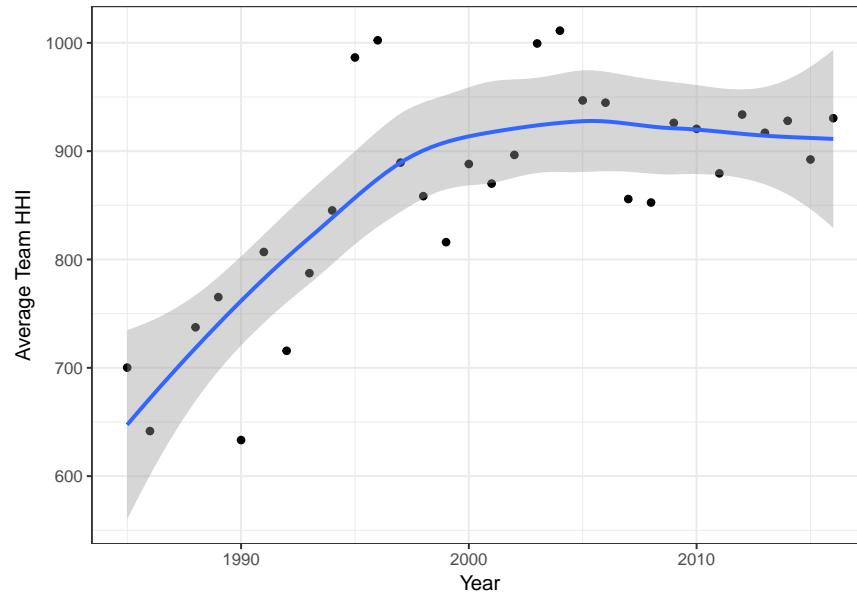
```
underlay_data <- filter(team_salary_vs_time, !flag)
overlay_data <- filter(team_salary_vs_time, flag)
```

```
ggplot() +
  geom_point(data = underlay_data, aes(x = yearID, y = avg, color = teamID)) +
  geom_line(data = underlay_data, aes(x = yearID, y = avg, color = teamID)) +
  geom_point(data = overlay_data, aes(x = yearID, y = avg, color = teamID)) +
  geom_line(data = overlay_data, aes(x = yearID, y = avg, color = teamID)) +
  scale_y_continuous(labels = comma) +
  scale_color_manual(values = team_colors) +
  labs(color = 'Team') +
  xlab('Year') +
  ylab('Average Team Salary')
```

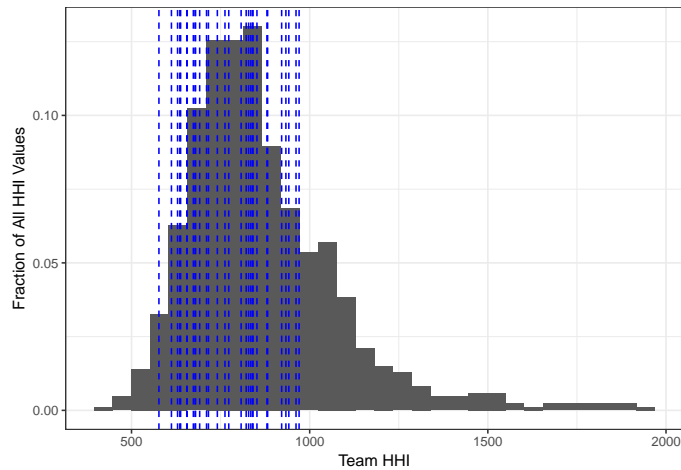


```
hhi_vs_time <- teams %>%
  group_by(yearID) %>%
  summarize(avg = mean(HHI))
```

```
ggplot(data = hhi_vs_time) +
  geom_point(aes(x = yearID, y = avg)) +
  geom_smooth(aes(x = yearID, y = avg)) +
  xlab('Year') +
  ylab('Average Team HHI')
```



```
ggplot(data = filter(teams, mean(teams$HHI) - 5 * sd(teams$HHI) <= HHI & HHI <= mean(teams$HHI) + 5 * sd(teams$HHI))) +
  geom_histogram(aes(x = HHI, y = (..count..) / sum(..count..))) +
  geom_vline(data = filter(teams, WSWin), aes(xintercept = HHI), color = 'blue', linetype = 'dashed') +
  xlab('Team HHI') +
  ylab('Fraction of All HHI Values')
```



```
year_to_period <- function(year) {
  if (year <= 1992) {
    return('Pre 1993')
  }
  else if (1993 <= year & year <= 1997) {
    return('Between 1993 and 1997')
  }
  else {
    return('Post 1997')
  }
}
```

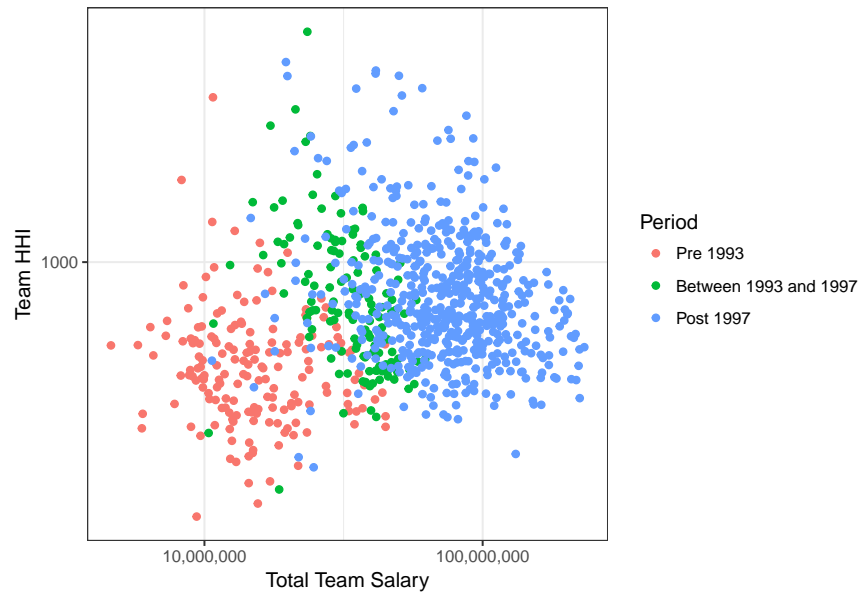
```
hhi_vs_total_salary <- teams %>%
```

```

mutate(period = year_to_period(yearID))
hhi_vs_total_salary$period <- factor(hhi_vs_total_salary$period, levels = c('Pre 1993', 'Between 1993 and 1997', 'Post 1997'))

ggplot(data = hhi_vs_total_salary) +
  geom_point(aes(x = totalSalaryMil * 1000000, y = HHI, color = period)) +
  scale_x_log10(labels = comma) +
  scale_y_log10() +
  labs(color = 'Period') +
  xlab('Total Team Salary') +
  ylab('Team HHI')

```



```

teams_new <- teams %>%
  filter(1999 <= yearID & yearID <= 2016) %>%
  mutate(normalizedYear = yearID - 1999)

```

```

salaries_new <- salaries %>%
  filter(1999 <= yearID & yearID <= 2016)

```

```

linear_fixed_new <- lm(formula = winPercentage ~ totalSalaryMil + HHI +
  normalizedYear + teamID + 0,
  data = teams_new)
summary(linear_fixed_new)$coefficients[1:3,]

```

```

##              Estimate Std. Error  t value    Pr(>|t|)
## totalSalaryMil  0.49949048 0.13266465  3.765061 0.0001868427
## HHI            -0.05358433 0.01442365 -3.715033 0.0002267173
## normalizedYear -2.00506563 0.73422847 -2.730847 0.0065462996

```

```

linear_random_new <- lm(formula = winPercentage ~ totalSalaryMil + HHI + normalizedYear,
  data = teams_new)
summary(linear_random_new)$coefficients[1:4,]

```

```

##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept)  503.61973192 15.73186679 32.012713 7.433408e-125
## totalSalaryMil  0.70690261 0.08558449  8.259705 1.226558e-15

```



```
## HHI          -0.04403327  0.01403935 -3.136419  1.807431e-03
## normalizedYear -2.75511624  0.63636357 -4.329469  1.794894e-05

#Gini coef: https://en.wikipedia.org/wiki/Gini\_coefficient

gini_coef <- salaries %>%
  group_by(yearID, teamID) %>%
  summarize(gini_coef = Gini(salary))

teams <- teams %>%
  inner_join(gini_coef)

## Joining, by = c("yearID", "teamID")
#Compare Gini index to HHI for old teams
gini_old <- teams %>%
  filter(1985 <= yearID & yearID <= 1998) %>%
  mutate(normalizedYear = yearID - 1985)

gini_fixed_old <- lm(formula = winPercentage ~ totalSalaryMil + gini_coef +
                     normalizedYear + teamID + 0,
                     data = gini_old)
summary(gini_fixed_old)$coefficients[1:3,]

##              Estimate Std. Error  t value    Pr(>|t|)
## totalSalaryMil    2.312489   0.4521101   5.114879 5.559085e-07
## gini_coef        -129.827687  58.3962938  -2.223218 2.693371e-02
## normalizedYear    -3.793613   1.8442237  -2.057025 4.053398e-02

gini_random_old <- lm(formula = winPercentage ~ totalSalaryMil + gini_coef + normalizedYear,
                     data = gini_old)
summary(gini_random_old)$coefficients[1:4,]

##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept)      539.115522  24.3427437  22.146868 1.404711e-67
## totalSalaryMil    2.495248   0.3775998   6.608182 1.532621e-10
## gini_coef        -133.160525  54.5170461  -2.442548 1.510025e-02
## normalizedYear    -4.559438   1.6394568  -2.781066 5.724596e-03

#Compare Gini coef to HHI results for new teams
gini_new <- teams %>%
  filter(1999 <= yearID & yearID <= 2016) %>%
  mutate(normalizedYear = yearID - 1999)

gini_fixed_new <- lm(formula = winPercentage ~ totalSalaryMil + gini_coef +
                     normalizedYear + teamID + 0,
                     data = gini_new)
summary(gini_fixed_new)$coefficients[1:3,]

##              Estimate Std. Error  t value    Pr(>|t|)
## totalSalaryMil    0.6624573   0.1276041   5.191504 3.069387e-07
## gini_coef        -218.5193970  58.0041403  -3.767307 1.852177e-04
## normalizedYear    -2.7481628   0.7167617  -3.834137 1.425403e-04

gini_random_new <- lm(formula = winPercentage ~ totalSalaryMil + gini_coef + normalizedYear,
                     data = gini_new)
summary(gini_random_new)$coefficients[1:4,]
```

```
##               Estimate Std. Error   t value    Pr(>|t|)
## (Intercept)    550.1960765 32.40207090 16.980275 9.795202e-52
## totalSalaryMil    0.8147846  0.08086385 10.076006 6.326317e-22
## gini_coef     -158.9308612 55.31835762 -2.873022 4.232464e-03
## normalizedYear   -3.2653945  0.62553010 -5.220204 2.591533e-07

#Atkinson index: https://en.wikipedia.org/wiki/Atkinson\_index
atkinson <- salaries %>%
  group_by(yearID, teamID) %>%
  summarize(atk = Atkinson(salary))

teams <- teams %>%
  inner_join(atkinson)

## Joining, by = c("yearID", "teamID")
#Compare Atkinson index to HHI for old teams
atk_old <- teams %>%
  filter(1985 <= yearID & yearID <= 1998) %>%
  mutate(normalizedYear = yearID - 1985)

atk_fixed_old <- lm(formula = winPercentage ~ totalSalaryMil + atk +
                    normalizedYear + teamID + 0,
                    data = atk_old)
summary(atk_fixed_old)$coefficients[1:3,]

##               Estimate Std. Error   t value    Pr(>|t|)
## totalSalaryMil    2.376013  0.4504696  5.274524 2.527538e-07
## atk              -179.252498 65.9775739 -2.716870 6.966615e-03
## normalizedYear   -3.413854  1.8123452 -1.883667 6.056128e-02

atk_random_old <- lm(formula = winPercentage ~ totalSalaryMil + atk + normalizedYear,
                    data = atk_old)
summary(atk_random_old)$coefficients[1:4,]

##               Estimate Std. Error   t value    Pr(>|t|)
## (Intercept)    509.553006 11.7544621 43.349751 1.975329e-139
## totalSalaryMil    2.555104  0.3759984  6.795518 4.948012e-11
## atk             -181.545622 61.0995554 -2.971308 3.179705e-03
## normalizedYear   -4.188912  1.6071160 -2.606478 9.556548e-03

#Compare Atkinson coef to HHI results for new teams
atk_new <- teams %>%
  filter(1999 <= yearID & yearID <= 2016) %>%
  mutate(normalizedYear = yearID - 1999)

atk_fixed_new <- lm(formula = winPercentage ~ totalSalaryMil + atk +
                    normalizedYear + teamID + 0,
                    data = atk_new)
summary(atk_fixed_new)$coefficients[1:3,]

##               Estimate Std. Error   t value    Pr(>|t|)
## totalSalaryMil    0.6943805  0.1274283  5.449185 8.049458e-08
## atk              -267.2160401 60.5183651 -4.415454 1.242698e-05
## normalizedYear   -2.8583544  0.7142924 -4.001659 7.268982e-05
```

```
atk_random_new <- lm(formula = winPercentage ~ totalSalaryMil + atk + normalizedYear,
                     data = atk_new)
summary(atk_random_new)$coefficients[1:4,]
```

```
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)   512.1999193  17.31125767  29.587678 2.179218e-113
## totalSalaryMil    0.8479434   0.08182447  10.362956 5.387417e-23
## atk            -193.0102061  57.91250941  -3.332790 9.212586e-04
## normalizedYear   -3.3980916   0.62691299  -5.420356 9.131764e-08
```

```
teams <- teams %>%
  mutate(normalizedYear = yearID - 1985) %>%
  ungroup()

set.seed(42)
num_folds <- 5
num_rows <- nrow(teams)
ndx <- sample(1:num_rows, num_rows, replace=F)

train_data <- teams[ndx, ] %>%
  mutate(fold = (row_number() %% num_folds) + 1)

validate_err <- c()
train_err <- c()
for (f in 1:num_folds) {
  temp_train <- filter(train_data, fold != f)

  #Using fixed effects model because the data is randomly sorted
  model <- lm(formula = winPercentage ~ totalSalaryMil + HHI +
              normalizedYear + teamID + 0, data = temp_train)
  train_err[f] <- sqrt(mean( (predict (model, temp_train) - temp_train$winPercentage) ^2) )

  # evaluate on the validation data
  temp_validate <- filter(train_data, fold == f)
  validate_err[f] <- sqrt(mean( (predict (model, temp_validate) - temp_validate$winPercentage) ^2) )
}

# compute the average validation error across folds and the standard error on this estimate
avg_validate_err <- mean(validate_err)
se_validate_err <- sd(validate_err) / sqrt(num_folds)

avg_train_err <- mean(train_err)
se_train_err <- sd(train_err) / sqrt(num_folds)

time_train_data <- teams %>% filter(yearID < 2012)
time_validate_data <- teams %>% filter(yearID>2011)

#Using random effects model because we dont want to worry about team being good in 80s/90s and bad in 2000s
model <- lm(formula = winPercentage ~ totalSalaryMil + HHI +
            normalizedYear,
            data = time_train_data)

time_train_err <- sqrt(mean( (predict (model, time_train_data) - time_train_data$winPercentage) ^2) )
```

```
time_validate_err <- sqrt(mean( (predict (model, time_validate_data) - time_validate_data$winPercentage
```

Extension Notes

- note that minimum salary has increased over time: https://www.baseball-reference.com/bullpen/Minimum_salary
- Coefficients of different inequality indexes are the same sign and have the same predictive impact (magnitudes are different because HHI has a wider range than the other two which are between 0 and 1)
- Used fixed effects model when split data independent of time, ie teams from 1985 and 2016 are in the training set. Meanwhile used random effects model when split data based on time, ie teams from 1985-2011 were in training set and 2012-2016 were test set. The use of the random effects model for the latter was to account for the fact that teams might be dominant in early years but not so much in more recent years. Also handles the issue of teams not existing in training set but potentially in validation set.

Extension Analysis

- 6.4%-6.6% error in win percentage predictions for both kinds of regressions. This isn't amazing given that the range of reasonable values is roughly 40%-70% win percentage

Postface

The following is a list of all packages used to generate these results.

```
sessionInfo()
```

```
## R version 3.4.3 (2017-11-30)
## Platform: x86_64-apple-darwin15.6.0 (64-bit)
## Running under: macOS Sierra 10.12.6
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/3.4/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/3.4/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  base
##
## other attached packages:
## [1] bindrcpp_0.2      ineq_0.2-13      forcats_0.3.0    stringr_1.3.0
## [5] dplyr_0.7.4       purrr_0.2.4      readr_1.1.1      tidyr_0.8.0
## [9] tibble_1.4.2      ggplot2_2.2.1    tidyverse_1.2.1  scales_0.5.0
## [13] here_0.1
##
## loaded via a namespace (and not attached):
## [1] Rcpp_0.12.15      cellranger_1.1.0 compiler_3.4.3    pillar_1.2.1
## [5] plyr_1.8.4        bindr_0.1         methods_3.4.3     tools_3.4.3
## [9] digest_0.6.15     lubridate_1.7.3   jsonlite_1.5      gtable_0.2.0
```

```

## [13] evaluate_0.10.1  nlme_3.1-131      lattice_0.20-35    pkgconfig_2.0.1
## [17] rlang_0.2.0      psych_1.7.8        cli_1.0.0          rstudioapi_0.7
## [21] yaml_2.1.17      parallel_3.4.3     haven_1.1.1        xml2_1.2.0
## [25] httr_1.3.1       knitr_1.20         hms_0.4.1          rprojroot_1.3-2
## [29] grid_3.4.3       glue_1.2.0         R6_2.2.2           readxl_1.0.0
## [33] foreign_0.8-69   rmarkdown_1.9      modelr_0.1.1       reshape2_1.4.3
## [37] magrittr_1.5     backports_1.1.2    htmltools_0.3.6    rvest_0.3.2
## [41] assertthat_0.2.0 mnormt_1.5-5       colorspace_1.3-2   labeling_0.3
## [45] stringi_1.1.6    lazyeval_0.2.1     munsell_0.4.3      broom_0.4.3
## [49] crayon_1.3.4

```