

Adjusted Home Run Frequencies

Jason Osborne and Rich Levine

Adjusted hr frequencies, 2012-2021.

Let us see how the batter-pitcher combination (bpcombo) frequencies have varied over time

```
##          PIT_HAND_CD
## BAT_HAND_CD      L      R
##          L 0.08054097 0.34107305
##          R 0.19685160 0.38153439

## [1] "by Batter Hand (row sums)"

##          L      R
## 0.421614 0.578386

## [1] "by Pitcher Hand (col sums)"

##          L      R
## 0.2773926 0.7226074

## [1] "Conditionally on Batter Hand"

##          PIT_HAND_CD
## BAT_HAND_CD      L      R
##          L 0.1910301 0.8089699
##          R 0.3403464 0.6596536

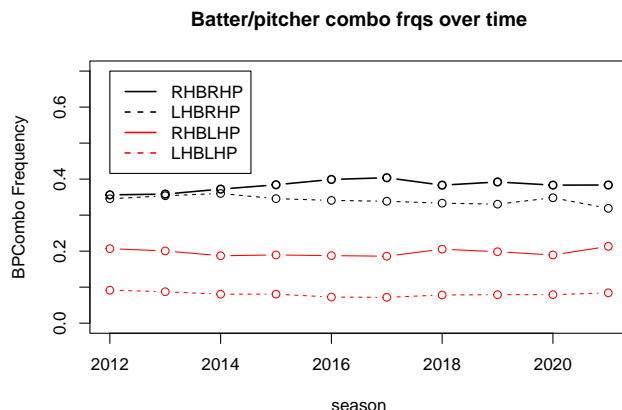
## [1] "Conditionally on Pitcher Hand"

##          PIT_HAND_CD
## BAT_HAND_CD      L      R
##          L 0.2903501 0.4720032
##          R 0.7096499 0.5279968

## [1] "All four relative freqs as a vector"

##  BAT_HAND_CD PIT_HAND_CD      Freq
## 1          L          L 0.08054097
## 2          R          L 0.19685160
## 3          L          R 0.34107305
## 4          R          R 0.38153439
```

These four bpcombo frequencies have changed little over time, though the preference for LHB when facing RHP may have decreased slightly.

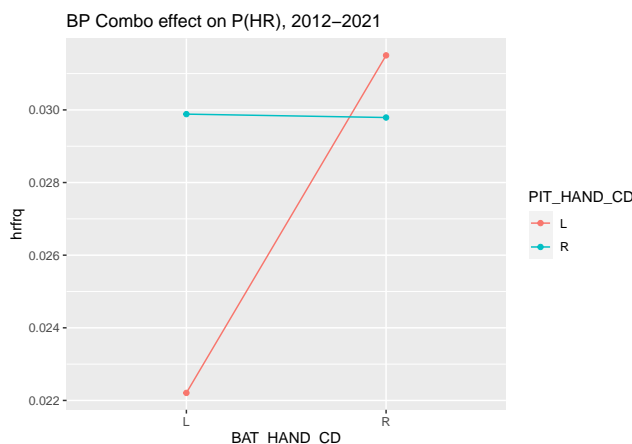


The effects of bpcombo on home run frequency can be investigated with an interaction plot, generated with data from 2012-2021:

```
allyrs.12vars %>% group_by(BAT_HAND_CD,PIT_HAND_CD) %>% summarize(hrfrq=mean(hr)) -> hrfrq.bp.era
hrfrq.bp.era
```

```
## # A tibble: 4 x 3
## # Groups:   BAT_HAND_CD [2]
##   BAT_HAND_CD PIT_HAND_CD hrfrq
##   <chr>        <chr>      <dbl>
## 1 L          L          0.0222
## 2 L          R          0.0299
## 3 R          L          0.0315
## 4 R          R          0.0298
```

```
ggplot(hrfrq.bp.era,aes(y=hrfrq,x=BAT_HAND_CD,color=PIT_HAND_CD)) +
  geom_line(aes(group=PIT_HAND_CD)) + geom_point() +
  ggtitle("BP Combo effect on P(HR), 2012-2021")
```



Looking over this 10 year period, it can be seen that home runs are least likely when a LHB is facing a LHP. Remarkably, the effect of the batter hand only appears to matter when facing lefties. Wow!

For a given park, the frequencies of the four combinations can vary dramatically from one season to the next, depending upon the personnel of the home team and with the unbalanced schedules of years past, upon the personnel of other teams in the division. In light of bpcombo effects, home run frequencies for each park can be adjusted to league-wide bpcombo frequencies simply by reweighting the four conditional home run rates

to the these frequencies.

```
allyrs.12vars %>% group_by(park,BAT_HAND_CD,PIT_HAND_CD) %>%
  summarize(hrfrq=mean(hr)) %>% pivot_wider(values_from=hrfrq,
                                             names_from=c("BAT_HAND_CD","PIT_HAND_CD")) ->
  hrsummary.wide
hrsummary.wide %>% mutate(adjhr=bpfrqs.era[1,1]*L_L + bpfrqs.era[1,2]*L_R +
                          bpfrqs.era[2,1]*R_L + bpfrqs.era[2,2]*R_R) ->
  hrsummary.wide
# unadjusted
allyrs.12vars %>% group_by(park) %>% summarize(obshr=mean(hr)) -> hrfrqs.bypark
hrsummary.wide %>% inner_join(hrfrqs.bypark) -> hrsummary.wide
```

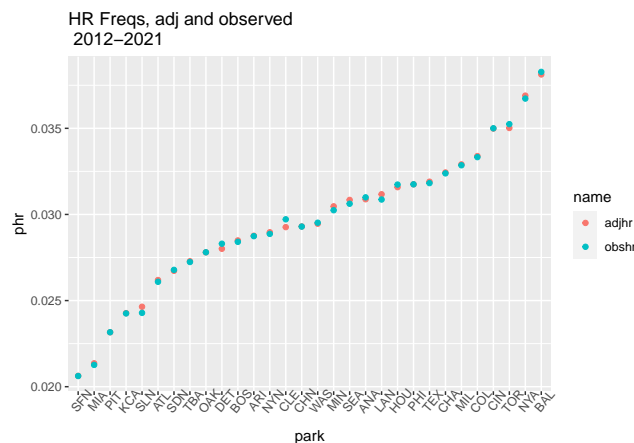
These adjusted frequencies can be plotted against park, along with the unadjusted frequencies. Further investigation of changes over time is warranted though, as a glance at 10-year averages still shows considerable variability in bpcombo frequencies across parks. It must be kept in mind that many players reside with the same team for long periods of time, so these 10 years are not at all independent. However, we average anyway ...

A technique worth mentioning in the construction of this plot is to achieve an ordering of parks on the horizontal axis according to either the observed or adjusted home run rate by so ordering the levels of park as a factor.

```
hrsummary.wide$park <- factor(hrsummary.wide$park,
                             levels=hrsummary.wide$park[order(hrsummary.wide$adjhr)])
hrsummary.wide %>% pivot_longer(cols=c("adjhr","obshr"),values_to="phr") ->
  hrsummary.tall
```

Now ggplot can be used ...

```
ggplot(hrsummary.tall) + geom_point(aes(y=phr,x=park,color=name)) +
  theme(axis.text.x=element_text(angle=50)) +
  ggtitle("HR Freqs, adj and observed \n 2012-2021")
```



Ok, let us consider those teams for which the observed and adjusted hr freqs were different.

```
hrsummary.wide %>% mutate(adjmnt=adjhr-obshr) %>%
  arrange(abs(adjmnt)) -> hrsummary.wide; hrsummary.wide %>% tail
```

```
## # A tibble: 6 x 8
## # Groups:   park [6]
##   park    L_L    L_R    R_L    R_R adjhr obshr adjmnt
##   <fct> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 TOR   0.0252 0.0330 0.0359 0.0384 0.0350 0.0352 -0.000223
## 2 MIN   0.0187 0.0275 0.0352 0.0332 0.0305 0.0302  0.000225
## 3 DET   0.0218 0.0261 0.0342 0.0278 0.0280 0.0283 -0.000294
## 4 LAN   0.0266 0.0340 0.0294 0.0305 0.0312 0.0309  0.000307
## 5 SLN   0.0184 0.0243 0.0312 0.0229 0.0246 0.0243  0.000357
## 6 CLE   0.0184 0.0322 0.0290 0.0291 0.0293 0.0297 -0.000453
```

These differences between observed relative frequencies are small, but the number of plate appearances is large:

```
allyrs.12vars %>% group_by(park) %>% summarize(pa=n(),hr=sum(hr)) %>%
  inner_join(hrsummary.wide) %>% mutate(hrdiff=adjmnt*pa) %>%
  arrange(abs(hrdiff)) -> hrsummary.wide ; hrsummary.wide %>% tail
```

```
## # A tibble: 6 x 11
##   park    pa    hr    L_L    L_R    R_L    R_R adjhr obshr adjmnt hrdiff
##   <chr> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 TOR   57373  2022 0.0252 0.0330 0.0359 0.0384 0.0350 0.0352 -0.000223 -12.8
## 2 MIN   58421  1767 0.0187 0.0275 0.0352 0.0332 0.0305 0.0302  0.000225  13.1
## 3 DET   57347  1623 0.0218 0.0261 0.0342 0.0278 0.0280 0.0283 -0.000294 -16.9
## 4 LAN   56243  1736 0.0266 0.0340 0.0294 0.0305 0.0312 0.0309  0.000307  17.2
## 5 SLN   56990  1384 0.0184 0.0243 0.0312 0.0229 0.0246 0.0243  0.000357  20.4
## 6 CLE   57104  1697 0.0184 0.0322 0.0290 0.0291 0.0293 0.0297 -0.000453 -25.9
```

Each of the six teams that have the largest absolute adjustment have an extreme value either for proportion of LHB or proportion of LHP. First the parks that host the *fewest* plate appearances by LHB:

```
allyrs.12vars %>% select(park,BAT_HAND_CD) %>% table %>% prop.table(margin="park") %>%
  as.data.frame() %>% pivot_wider(values_from=Freq,names_from=BAT_HAND_CD) %>%
  arrange(L) -> BHbyPark
BHbyPark %>% head
```

```
## # A tibble: 6 x 3
##   park    L    R
##   <fct> <dbl> <dbl>
## 1 ANA   0.380 0.620
## 2 DET   0.382 0.618
## 3 TOR   0.383 0.617
## 4 CHA   0.390 0.610
## 5 HOU   0.396 0.604
## 6 MIA   0.397 0.603
```

Now for the *most* plate appearances by LHB

```
BHbyPark %>% tail
```

```
## # A tibble: 6 x 3
##   park      L      R
##   <fct> <dbl> <dbl>
## 1 SFN    0.444 0.556
## 2 PHI    0.455 0.545
## 3 NYN    0.457 0.543
## 4 MIN    0.458 0.542
## 5 SEA    0.468 0.532
## 6 CLE    0.519 0.481
```

Note that DET and TOR both see large downward adjustment and host the 2nd and 3rd lowest LHB frequencies at 38%. CLE and MIN see large upward adjustment and host the most and third most LHB, respectively (52% and 46% LHB!).

The other two teams for which adjustments are largest have extremes for plate appearances involving LHP:

```
all yrs.12vars %>% select(park,PIT_HAND_CD) %>% table %>% prop.table(margin="park") %>%
  as.data.frame() %>% pivot_wider(values_from=Freq,names_from=PIT_HAND_CD) %>%
  arrange(L) -> PHbyPark
PHbyPark %>% head
```

```
## # A tibble: 6 x 3
##   park      L      R
##   <fct> <dbl> <dbl>
## 1 MIL    0.206 0.794
## 2 SLN    0.206 0.794
## 3 CIN    0.206 0.794
## 4 CLE    0.215 0.785
## 5 MIA    0.243 0.757
## 6 NYN    0.243 0.757
```

```
PHbyPark %>% tail
```

```
## # A tibble: 6 x 3
##   park      L      R
##   <fct> <dbl> <dbl>
## 1 SFN    0.320 0.680
## 2 CHA    0.323 0.677
## 3 TEX    0.326 0.674
## 4 BOS    0.327 0.673
## 5 SEA    0.339 0.661
## 6 LAN    0.362 0.638
```

Dodger Stadium (LAN) has seen the greatest number of plate appearances with a LHP (36%) while Busch (SLN) has seen the second fewest (21%.) The variation in frequency of LHB across parks (38% for ANA up to 52% for CLE) and LHP (21% for MIL up to 36% for LAN) is remarkable. Lineups and rotations are perhaps more stable than one might think given all the personnel changes by high-profile free agents.

HR v park plots separated by batterhand

The adjusted HR Frequency for LHB is the weighted average of observed HR frequencies against LHP and RHP, with weights given by the conditionals from bpcombos on page 1

```
allyrs.12vars %>% select(BAT_HAND_CD,PIT_HAND_CD) %>% table %>%
  prop.table(margin="BAT_HAND_CD") -> bhtable

# compute observed hr frqs by park
allyrs.12vars %>% group_by(park,BAT_HAND_CD,PIT_HAND_CD) %>%
  summarize(hrfrq=mean(hr)) %>%
  pivot_wider(values_from=hrfrq,
              names_from=c("BAT_HAND_CD","PIT_HAND_CD")) ->
  hrsummary.wide
# compute weighted HR freq for LHB and for RHB
hrsummary.wide %>%
  mutate(adjhrLHB=bhtable[1,1]*L_L+bhtable[1,2]*L_R,
         adjhrRHB=bhtable[2,1]*R_L+bhtable[2,2]*R_R) ->
hrsummary.wide

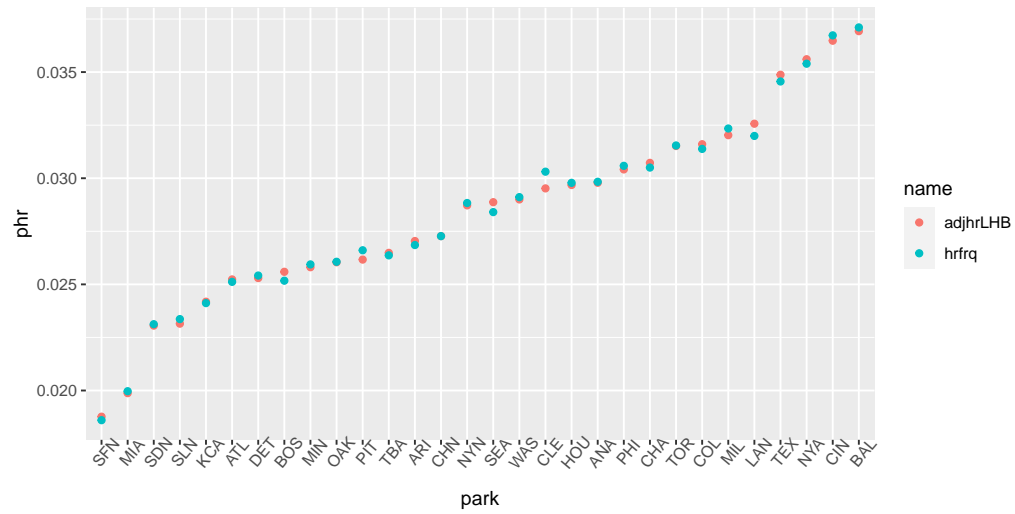
# hrfrs by hand not adjusted for pitcher hand
allyrs.12vars %>% group_by(park,BAT_HAND_CD) %>%
  summarize(hrfrq=mean(hr)) %>%
  inner_join(hrsummary.wide) -> hrsummary.wide

# make LHB tall
hrsummary.wide %>% filter(BAT_HAND_CD=="L") -> hrsummary.wide.LHB
hrsummary.wide.LHB$park <- factor(hrsummary.wide.LHB$park,
  levels=hrsummary.wide.LHB$park[order(hrsummary.wide.LHB$adjhrLHB)])

hrsummary.wide.LHB %>%
  pivot_longer(cols=c("adjhrLHB","hrfrq"),values_to="phr") ->
  hrsummary.tall.LHB

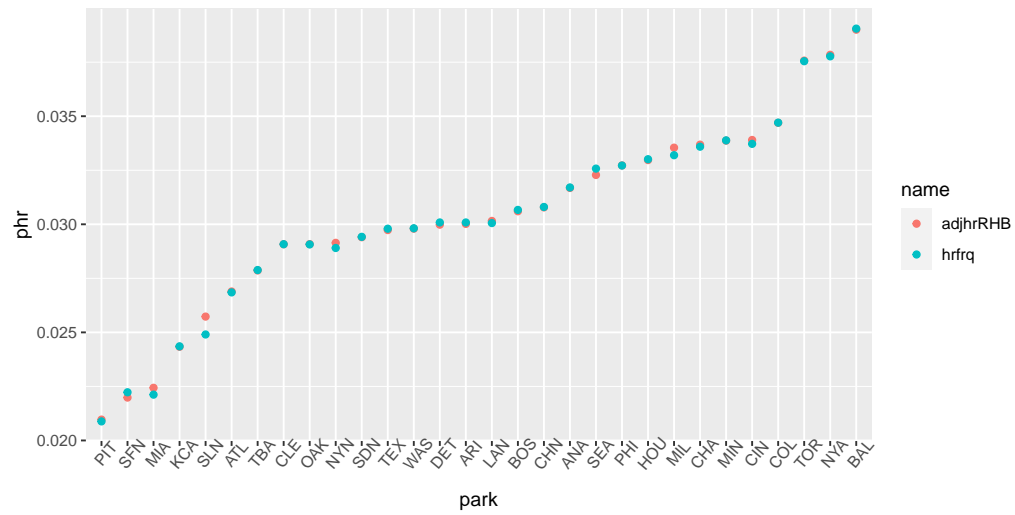
# plot LHB
ggplot(hrsummary.tall.LHB) + geom_point(aes(y=phr,x=park,color=name))+
  theme(axis.text.x=element_text(angle=50)) +
  ggtitle("HR Freqs for LHB, adjusted and observed \n 2012-2021")
```

HR Freqs for LHB, adjusted and observed
2012–2021



Similarly for RHB,

HR Freqs for RHB, adjusted and observed
2012–2021

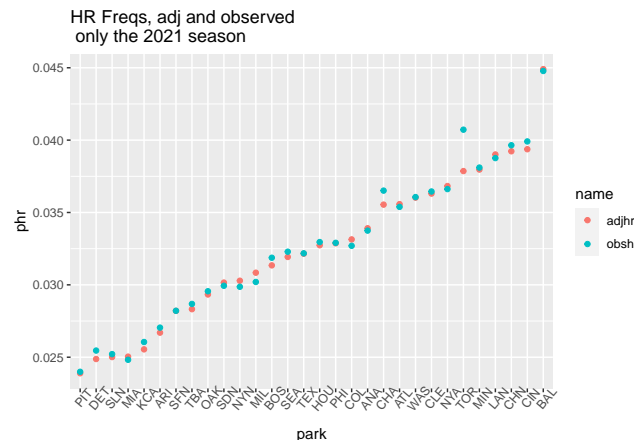


To investigate variability of adjusted and unadjusted HR freqs and also variability among rankings, we will here obtain graphs for the 2021 season alone. Firstly, HR freqs averaged over BH, to be followed by separate plots for LHB and RHB.

```
bpfrqs.2021 <- bpfrq.byyear[, ,10] # computed earlier
allyrs.12vars %>% filter(season==2021) %>%
  group_by(park,BAT_HAND_CD,PIT_HAND_CD) %>%
  summarize(hrfrq=mean(hr)) %>% pivot_wider(values_from=hrfrq,
                                             names_from=c("BAT_HAND_CD","PIT_HAND_CD")) ->
  hrsummary.wide.2021
hrsummary.wide.2021 %>% mutate(adjhr=bpfrqs.2021[1,1]*L_L + bpfrqs.2021[1,2]*L_R +
                              bpfrqs.2021[2,1]*R_L + bpfrqs.2021[2,2]*R_R) ->
  hrsummary.wide.2021
# unadjusted
allyrs.12vars %>% filter(season==2021) %>% group_by(park) %>% summarize(obshr=mean(hr)) -> hrfrqs.bypark.2021
hrsummary.wide.2021 %>% inner_join(hrfrqs.bypark.2021) -> hrsummary.wide.2021
hrsummary.wide.2021$park <- factor(hrsummary.wide.2021$park,
                                  levels=hrsummary.wide.2021$park[order(hrsummary.wide.2021$adjhr)])
hrsummary.wide.2021 %>% pivot_longer(cols=c("adjhr","obshr"),values_to="phr") ->
  hrsummary.tall.2021
```

Now ggplot can be used ...

```
ggplot(hrsummary.tall.2021) + geom_point(aes(y=phr,x=park,color=name)) +
  ggtitle("HR Freqs, adj and observed \n only the 2021 season") +
  theme(axis.text.x=element_text(angle=50))
```

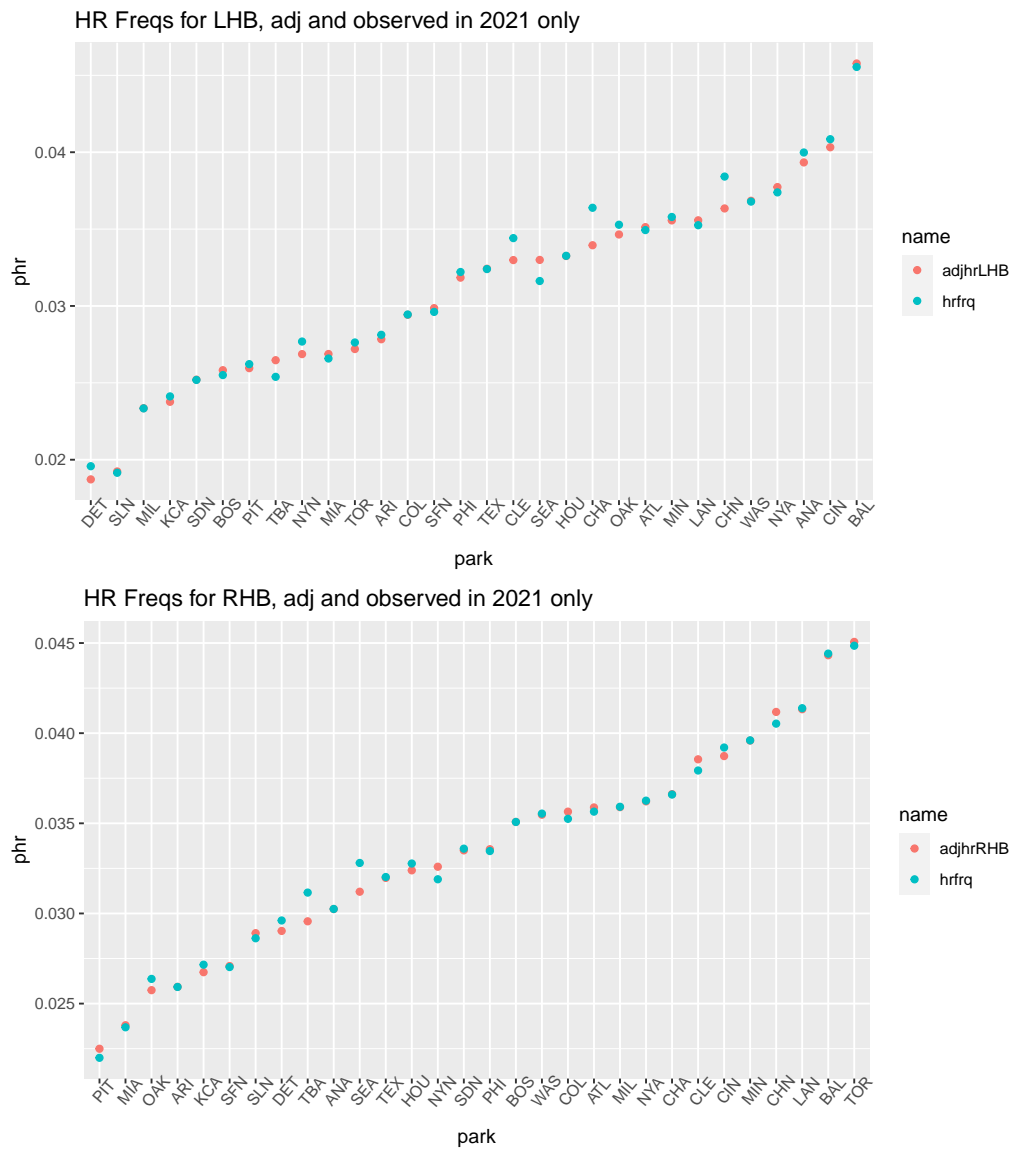


The data:

```
hrsummary.tall.2021 %>% print(n=5)

## # A tibble: 60 x 7
## # Groups:   park [30]
##   park    L_L    L_R    R_L    R_R name    phr
##   <fct> <dbl> <dbl> <dbl> <dbl> <chr> <dbl>
## 1 ANA    0.0465 0.0374 0.0323 0.0291 adjhr  0.0339
## 2 ANA    0.0465 0.0374 0.0323 0.0291 obshr  0.0337
## 3 ARI    0.0187 0.0303 0.0282 0.0246 adjhr  0.0267
## 4 ARI    0.0187 0.0303 0.0282 0.0246 obshr  0.0270
## 5 ATL    0.0302 0.0364 0.0312 0.0385 adjhr  0.0356
## # i 55 more rows
```


Now for 2021 plots specific to batter hand. Code available upon request!



Compute ranks for comparison of observed v fitted for LHB and for RHB

park	hrfrq	L_L	L_R	R_L	R_R	adjhrLHB	adjhrRHB	rLHBobs	rLHBfit
DET	0.020	0.005	0.022	0.033	0.027	0.019	0.029	2	1
SLN	0.019	0.021	0.019	0.032	0.027	0.019	0.029	1	2
MIL	0.023	0.024	0.023	0.041	0.033	0.023	0.036	3	3
KCA	0.024	0.012	0.027	0.034	0.023	0.024	0.027	4	4
SDN	0.025	0.021	0.026	0.032	0.034	0.025	0.033	5	5
BOS	0.026	0.013	0.029	0.035	0.035	0.026	0.035	7	6
PIT	0.026	0.019	0.028	0.030	0.018	0.026	0.022	8	7
TBA	0.025	0.015	0.029	0.038	0.025	0.026	0.030	6	8
NYN	0.028	0.008	0.032	0.040	0.028	0.027	0.033	11	9
MIA	0.027	0.015	0.030	0.027	0.022	0.027	0.024	9	10
TOR	0.028	0.032	0.026	0.042	0.047	0.027	0.045	10	11
ARI	0.028	0.019	0.030	0.028	0.025	0.028	0.026	12	12
COL	0.029	0.025	0.031	0.039	0.034	0.029	0.036	13	13
SFN	0.030	0.020	0.032	0.026	0.028	0.030	0.027	14	14
PHI	0.032	0.020	0.035	0.035	0.033	0.032	0.034	16	15
TEX	0.032	0.031	0.033	0.028	0.034	0.032	0.032	17	16
CLE	0.034	0.021	0.036	0.047	0.034	0.033	0.039	19	17
SEA	0.032	0.022	0.036	0.045	0.023	0.033	0.031	15	18
HOU	0.033	0.043	0.031	0.029	0.034	0.033	0.032	18	19
CHA	0.036	0.010	0.040	0.036	0.037	0.034	0.037	24	20
OAK	0.035	0.045	0.032	0.031	0.023	0.035	0.026	22	21
ATL	0.035	0.030	0.036	0.031	0.038	0.035	0.036	20	22
MIN	0.036	0.019	0.040	0.041	0.039	0.036	0.040	23	23
LAN	0.035	0.032	0.037	0.047	0.038	0.036	0.041	21	24
CHN	0.038	0.015	0.042	0.046	0.039	0.036	0.041	27	25
WAS	0.037	0.022	0.041	0.038	0.034	0.037	0.035	25	26
NYA	0.037	0.033	0.039	0.037	0.036	0.038	0.036	26	27
ANA	0.040	0.047	0.037	0.032	0.029	0.039	0.030	28	28
CIN	0.041	0.027	0.044	0.035	0.041	0.040	0.039	29	29
BAL	0.046	0.034	0.049	0.046	0.043	0.046	0.044	30	30

Some observations upon inspection of these plots and ranking

* Coors Field is middle of the road (#13) for LHB! * Yankee Stadium is indeed homer-friendly to LHB (), even more so if reweighting observed home run rates to league-wide matchup frequencies. * Jacobs Field (CLE) and T-Mobile (SEA) swith places in the ranking after adjustment, with CLE seeing a higher rate of observed HRs by LHB than would be expected for league-wide matchup frequencies * Both parks in Chicago saw more HRs by LHB than would be expected with league wide matchup frequencies

```
## # A tibble: 30 x 7
## # Groups:   park [30]
##   park BAT_HAND_CD hrfrq adjhrLHB adjhrRHB r_RHB_obs r_RHB_fit
##   <fct> <chr>      <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
## 1 TOR    R          0.0449  0.0272  0.0451      1      1
## 2 BAL    R          0.0444  0.0458  0.0443      2      2
## 3 LAN    R          0.0414  0.0356  0.0413      3      3
## 4 CHN    R          0.0405  0.0363  0.0412      4      4
## 5 MIN    R          0.0396  0.0356  0.0396      5      5
## 6 CIN    R          0.0392  0.0403  0.0387      6      6
## 7 CLE    R          0.0379  0.0330  0.0385      7      7
## 8 CHA    R          0.0366  0.0339  0.0366      8      8
```

##	9	NYA	R	0.0363	0.0377	0.0362	9	9
##	10	MIL	R	0.0359	0.0234	0.0359	10	10
##	11	ATL	R	0.0356	0.0351	0.0359	11	11
##	12	COL	R	0.0352	0.0294	0.0356	13	12
##	13	WAS	R	0.0355	0.0368	0.0355	12	13
##	14	BOS	R	0.0351	0.0258	0.0351	14	14
##	15	PHI	R	0.0335	0.0318	0.0336	16	15
##	16	SDN	R	0.0336	0.0252	0.0335	15	16
##	17	NYN	R	0.0319	0.0269	0.0326	20	17
##	18	HOU	R	0.0328	0.0332	0.0324	18	18
##	19	TEX	R	0.0320	0.0324	0.0320	19	19
##	20	SEA	R	0.0328	0.0330	0.0312	17	20
##	21	ANA	R	0.0302	0.0393	0.0302	22	21
##	22	TBA	R	0.0312	0.0265	0.0296	21	22
##	23	DET	R	0.0296	0.0187	0.0290	23	23
##	24	SLN	R	0.0286	0.0192	0.0289	24	24
##	25	SFN	R	0.0270	0.0299	0.0271	26	25
##	26	KCA	R	0.0272	0.0238	0.0267	25	26
##	27	ARI	R	0.0259	0.0278	0.0259	28	27
##	28	OAK	R	0.0264	0.0346	0.0257	27	28
##	29	MIA	R	0.0237	0.0269	0.0238	29	29
##	30	PIT	R	0.0220	0.0260	0.0225	30	30

I'd bet that the agreement between 2021 adjusted HRs and *era* observed HRs is greater than agreement between 2021 observed HRs and *era* observed HRs.

```
library(knitr)
knitr::knit_exit()
```