

# The Four Factors: Statistics that Make College Basketball Teams Win

Jake Singleton

11/12/2020

## Motivation: Dean Oliver and the NBA

In the early 2000s, statistician Dean Oliver published his findings on the “four-factor model,” a system used to evaluate the performance of NBA teams. He found that the four factors do a very nice job of revealing teams’ strengths and weaknesses in addition to being good predictors of wins. They are as follows:

1. Effective Field Goal Percentage (EFG)
2. Turnovers Committed per Possession (TPP)
3. Offensive Rebounding Percentage (ORP)
4. Free Throw Rate (FTR)

EFG is calculated as  $(\text{all field goals made} + 0.5(\text{3-point field goals made})) / (\text{all field goal attempts})$ . The great thing about EFG is that it weights 3-pointers appropriately, as 3-pointers are 50% more valuable than 2-pointers. Thus, it captures a team’s shooting ability far better than standard FG%.

TPP is simply  $(\text{turnovers committed}) / (\text{possessions})$ . Smaller TPP is better, as it means a team gives up the ball less often.

ORP is  $(\text{offensive rebounds}) / (\text{offensive rebounds} + \text{opponent’s defensive rebounds})$ . Obviously, high ORP is good and means the team gets second chances at making baskets.

FTR is  $(\text{free throw attempts}) / (\text{field goal attempts})$  and measures a team’s propensity to get to the free throw line (a good thing)!

Note that all of the numbers needed to calculate the four factors can be found in the box score. (Kind of... possessions are not in the box score, but analysts have come up with good formulas that estimate the number of possessions in a game very closely.) Not only does this makes our task easier, but it’s also interesting that a team’s win total can be explained so well by simple box score statistics (obviously these are modifications of box score stats, but they’re nothing fancy). While Oliver called these statistics the four factors, he calculated them from a defensive point of view as well, yielding eight factors. The defensive factors are:

1. Opponent’s Effective Field Goal Percentage (OPP\_EFG), i.e. the EFG a team yields to their opponents. Lower is better.
2. Defensive Turnovers Caused per Possession (DTPP). Higher is better.
3. Defensive Rebounding Percentage (DRP). Higher is better.
4. Opponent’s Free Throw Rate (OPP\_FTR). Lower is better.

These are computed analogously to their offensive counterparts.

So, how important are these factors? Oliver’s analysis led to the following rankings:

1. Shooting differential, EFG - OPP\_EFG (40% importance)
2. Turnover differential, TPP - DTPP (25%)
3. Rebounding differential, ORP - DPP (20%)
4. Free throw rate differential, FTR - OPP\_FTR (15%)

Note that other analysis online (e.g. [here](#)) and in Chapter 28 of Wayne Winston's *Mathletics* has shown Oliver's rankings overvalue free throw rate and rebounding and undervalue turnovers and shooting. Specifically, Winston found that EFG\_diff explains 71% of variance in NBA team wins, TO\_diff 15%, REB\_diff 6%, and FTR\_diff essentially 0%. I should also remind you that Oliver's work was done in the early-mid 2000s before the NBA's "3-point revolution" of the mid-late 2010s, and so I believe EFG\_diff is absolutely more valuable now than what it was.

In this analysis, I will analyze college basketball data from the 2019-2020 season to see how the college game differs from the pro game.

## The Data

There are a lot of Division 1 college basketball games in a season—too many, in fact, to go through every single one. So, I looked at the 10 “high-major” conferences: the Power 5 plus the Big East, Mountain West, Atlantic 10, American Athletic, and West Coast Conferences. From each conference, I took the champion, a middling team, and bottom of the barrel, yielding a total of 30 teams. These teams are: (Big12: Kansas, Texas Tech, KState), (WCC: Gonzaga, Pepperdine, Portland), (A10: Dayton, UMass, Fordham), (ACC: FSU, Syracuse, Wake Forest), (MountainW: SDSU, Colorado State, Wyoming), (SEC: Kentucky, A&M, Vandy), (Pac12: Oregon, Colorado, Washington), (Big10: Wisconsin, Penn State, Nebraska), (BigE: Creighton, Marquette, DePaul), and (AAC: Houston, SMU, Tulane). I employed this approach in order to include a wide variety of team talent. I limited the games scraped for each team to 30 games, and due to overlap, there are 832 games. All data comes from the 19-20 season.

## The Code

To scrape the relevant data, the key part was collecting all of the necessary ESPN game IDs. For example, California's first game of the season was against Pepperdine, and the box score is at this link: <https://www.espn.com/mens-college-basketball/matchup?gameId=401170553>. The last 9 digits of the link comprise the game ID, and I had to collect each one manually. If anyone knows a better way of getting these IDs, please let me know.

I'm not going to show all of the code since it's pretty long, but I will display all of the important outputs. You can find the R Markdown file in its full form in the code directory of this repository.

Team	FGM	3FGM	FGA	FTA	FTM	TO	OR	DR	OPP_FGM	OPP_3FGM
KU	23	4	50	26	16	28	10	30	23	8
DUKE	23	8	64	23	14	16	11	19	23	4
MONM	18	4	60	23	17	15	11	23	37	14
KU	37	14	66	31	24	7	11	32	18	4
ETSU	23	9	60	15	8	16	11	22	30	1
KU	30	1	54	18	14	15	5	27	23	9
KU	38	12	69	11	5	14	13	20	23	7
CHAM	23	7	50	14	10	27	6	22	38	12
KU	29	4	61	15	9	12	9	26	22	9
BYU	22	9	54	4	3	20	5	27	29	4

The above table shows a slice of the relevant box score statistics for each game that we'll use to calculate

the four factors. For example, the first two rows represent the Kansas vs. Duke game, one of the first of the season.

```
# Group data by team and aggregate (we use sum of course)
grouped = all_game_df %>%
  group_by(Team) %>%
  summarize(across(.cols = everything(), sum)) %>%
  filter(Team %in% teams_of_interest)

kable(grouped[1:5, 1:11], "simple") %>%
  kable_styling(font_size = 5)
```

Team	FGM	3FGM	FGA	FTA	FTM	TO	OR	DR	OPP_FGM	OPP_3FGM
COLO	730	221	1661	605	447	400	299	808	718	187
CREI	833	290	1776	525	391	322	238	800	788	229
CSU	789	241	1682	655	442	395	266	799	801	261
DAY	868	240	1664	576	414	367	230	816	702	187
DEP	770	168	1792	601	398	439	359	752	733	244

The code above aggregates the data for each team of interest by summing each column.

We see that this data frame includes each of the box score statistics necessary to calculate our factors. For example, Colorado made 730 field goals out of 1,661, and they gave up 718 successful field goals to their opponents. 221 of their 890 were 3-point field goals, and of their opponents' 718 makes, 187 were 3-pointers.

Team	Wins	EFG	TPP	ORP	FTR	OPP_EFG	DTPP	DRP	OPP_FTR
COLO	19	0.5060205	0.1972192	0.3029382	0.3642384	0.4712544	0.1850269	0.7651515	0.2694541
CREI	23	0.5506757	0.1539933	0.2361111	0.2956081	0.4828785	0.1754268	0.6999125	0.2306046
CSU	18	0.5407253	0.1881669	0.2646766	0.3894174	0.5198103	0.1793692	0.7624046	0.2421875
DAY	28	0.5937500	0.1786375	0.2575588	0.3461538	0.4657494	0.1882027	0.7351351	0.3167447
DEP	14	0.4765625	0.2054820	0.3146363	0.3353795	0.4956522	0.2228702	0.7008388	0.3553623
FOR	8	0.4489545	0.1907853	0.2086466	0.2607626	0.4703283	0.2114921	0.7342995	0.3364899
FSU	25	0.5192415	0.1903151	0.3349705	0.2950363	0.4681432	0.2400122	0.7000964	0.3483010
GONZ	28	0.5667929	0.1564844	0.3320274	0.3834505	0.4778498	0.1841731	0.7670251	0.2166397
HOU	22	0.4796574	0.1790741	0.3873484	0.3292291	0.4391646	0.1849033	0.7471060	0.4164649
KSU	9	0.4741538	0.2199956	0.2996183	0.3483077	0.5012937	0.2450883	0.6923879	0.3887451
KU	27	0.5406121	0.1905533	0.3285714	0.3513832	0.4395322	0.1957313	0.7362832	0.2338848
MARQ	18	0.5155896	0.1898605	0.3000941	0.3922902	0.4684611	0.1535202	0.7404903	0.3485590
MASS	13	0.4892254	0.1914545	0.2559923	0.2999418	0.5129260	0.2033031	0.6997027	0.3349001
NEB	6	0.4678629	0.1627741	0.2200489	0.2849491	0.5209596	0.1876514	0.6692573	0.2141414
ORE	23	0.5438845	0.1768460	0.3296482	0.2938845	0.4801865	0.2036620	0.6863680	0.2966200
PEPP	14	0.4966942	0.1761960	0.2661142	0.3256198	0.5269815	0.1904314	0.7210425	0.3395166
PORT	7	0.4839394	0.2001255	0.2037218	0.3315152	0.5144686	0.1864114	0.6953642	0.3038397
PSU	20	0.4919915	0.1567383	0.2894281	0.3064602	0.4667969	0.1845088	0.7253521	0.3203125
SDSU	28	0.5483967	0.1649777	0.2786017	0.2885986	0.4584644	0.2153056	0.7441176	0.2951542
SMU	19	0.5205158	0.1999960	0.3412698	0.3171161	0.4918746	0.1880449	0.7095865	0.2820662
SYR	16	0.4982984	0.1633413	0.2948244	0.3505389	0.4871134	0.2048096	0.6738739	0.3109966
TA&M	16	0.4640127	0.2159894	0.3132296	0.4108280	0.4915612	0.2256006	0.6819961	0.2603978
TTU	17	0.5132053	0.1995146	0.2786378	0.3451381	0.4624846	0.2343127	0.7019417	0.3579336
TULN	12	0.4780573	0.1710052	0.2229299	0.3417203	0.5269231	0.2322457	0.6878147	0.2958580
UK	24	0.5048164	0.1811274	0.3063830	0.4039735	0.4581197	0.1730189	0.7070347	0.3259259
VAN	9	0.4954901	0.1878194	0.2516370	0.4233313	0.5273048	0.1927910	0.7096774	0.3746330

Team	Wins	EFG	TPP	ORP	FTR	OPP_EFG	DTPP	DRP	OPP_FTR
WAKE	12	0.4831591	0.1976783	0.2847866	0.3983740	0.4961454	0.1535331	0.7373108	0.3458150
WASH	13	0.5054545	0.2163281	0.2736733	0.3848485	0.4418412	0.1923443	0.6803069	0.3224917
WISC	20	0.5063139	0.1556363	0.2378109	0.2633794	0.4670428	0.1696002	0.7410972	0.2577197
WYO	7	0.4821092	0.1864610	0.1453831	0.3038293	0.5122828	0.1893272	0.7117296	0.3786699

These are our four factors for each team (or eight, if you like). Note we have each team's number of wins now, which will come into play soon as our dependent variable.

Team	Wins	EFG	TPP	ORP	FTR	OPP_EFG	DTPP	DRP	OPP_FTR
DAY	28	1	11	21	13	7	17	10	15
GONZ	28	2	3	4	8	14	24	1	2
SDSU	28	4	7	17	27	5	7	5	10
KU	27	7	19	6	10	2	12	9	4
FSU	25	9	18	3	25	10	2	20	23
UK	24	15	13	9	3	4	27	17	18
CREI	23	3	1	25	24	16	26	21	3
ORE	23	5	10	5	26	15	10	26	12
HOU	22	24	12	1	18	1	22	4	30
PSU	20	19	4	14	21	8	23	12	16
WISC	20	12	2	24	29	9	28	6	6
COLO	19	13	22	10	9	13	21	2	8
SMU	19	8	25	2	20	19	18	16	9
CSU	18	6	16	20	6	26	25	3	5
MARQ	18	10	17	11	5	11	30	7	24
TTU	17	11	24	16	14	6	3	18	26
SYR	16	16	6	13	11	17	9	29	14
TA&M	16	29	28	8	2	18	5	27	7
DEP	14	26	27	7	16	20	6	19	25
PEPP	14	17	9	19	19	29	15	13	21
MASS	13	20	21	22	23	24	11	22	19
WASH	13	14	29	18	7	3	14	28	17
TULN	12	25	8	26	15	28	4	25	11
WAKE	12	22	23	15	4	21	29	8	22
KSU	9	27	30	12	12	22	1	24	29
VAN	9	18	15	23	1	30	13	15	27
FOR	8	30	20	28	30	12	8	11	20
PORT	7	21	26	29	17	25	20	23	13
WYO	7	23	14	30	22	23	16	14	28
NEB	6	28	5	27	28	27	19	30	1

This data frame is a table of rankings. For example, Dayton was 1st among our 30 teams in EFG, 12th in TPP, 22nd in ORP, 16th in FTR, etc.

## The Analysis

Interestingly, we see the top 3 teams in wins, Dayton, Gonzaga, and SDSU, are 1st, 2nd, and 4th respectively in EFG ranking. Moreover, 7 of the top 10 teams in wins are top 10 in EFG, and 6 of the top 10 rank inside the top 10 in OPP\_EFG. This supports the NBA analysis that has found shooting is the most important factor. If we look at the bottom six teams or so, not only do we see low ranks for EFG, but we also see the

same for ORP. Note that Fordham, Portland, Wyoming, and Nebraska all shot poorly and didn't get second chances either. Fordham is also interesting because they were last in EFG and FTR, yet they played good defense, ranking 9th and 8th in OPP\_EFG and DTPP respectively. As for a couple teams in the middle, SMU and Colorado State were good shooting teams (8th and 6th respectively), but poor defending shooters (19th and 26th respectively). Further, they were in the bottom half in turnovers committed on offense (25th and 16th). Another team I want to point out is Washington, the most unlucky team in the nation last year. They were 14th in EFG (not bad), 7th in FTR, and 3rd(!) in OPP\_EFG, but they finished last in the PAC-12 and won only 13 games. Note that they were 29th in TPP, 18th in ORP, and 28th in DRP, though, meaning they turned the ball over a lot and were not a good rebounding team.

Now, we will dig into some of the good statistics stuff, emulating Wayne Winston's analysis from his fantastic book *Mathletics*. We will run a regression to see which factors most influence wins. Our explanatory variables are:

1. Shooting differential: EFG - OPP\_EFG
2. Turnover differential: TPP - DTPP
3. Rebounding differential: ORP - DRP
4. Free throw differential: FTR - OPP\_FTR

Our response variable is wins.

```
# Prepare to run regression by making appropriate variables
deviation_tab = four_factors_df %>%
  mutate(EFG_diff = EFG - OPP_EFG,
         TO_diff = TPP - DTPP,
         REB_diff = ORP - DRP,
         FTR_diff = FTR - OPP_FTR) %>%
  select(Wins, EFG_diff, TO_diff, REB_diff, FTR_diff)

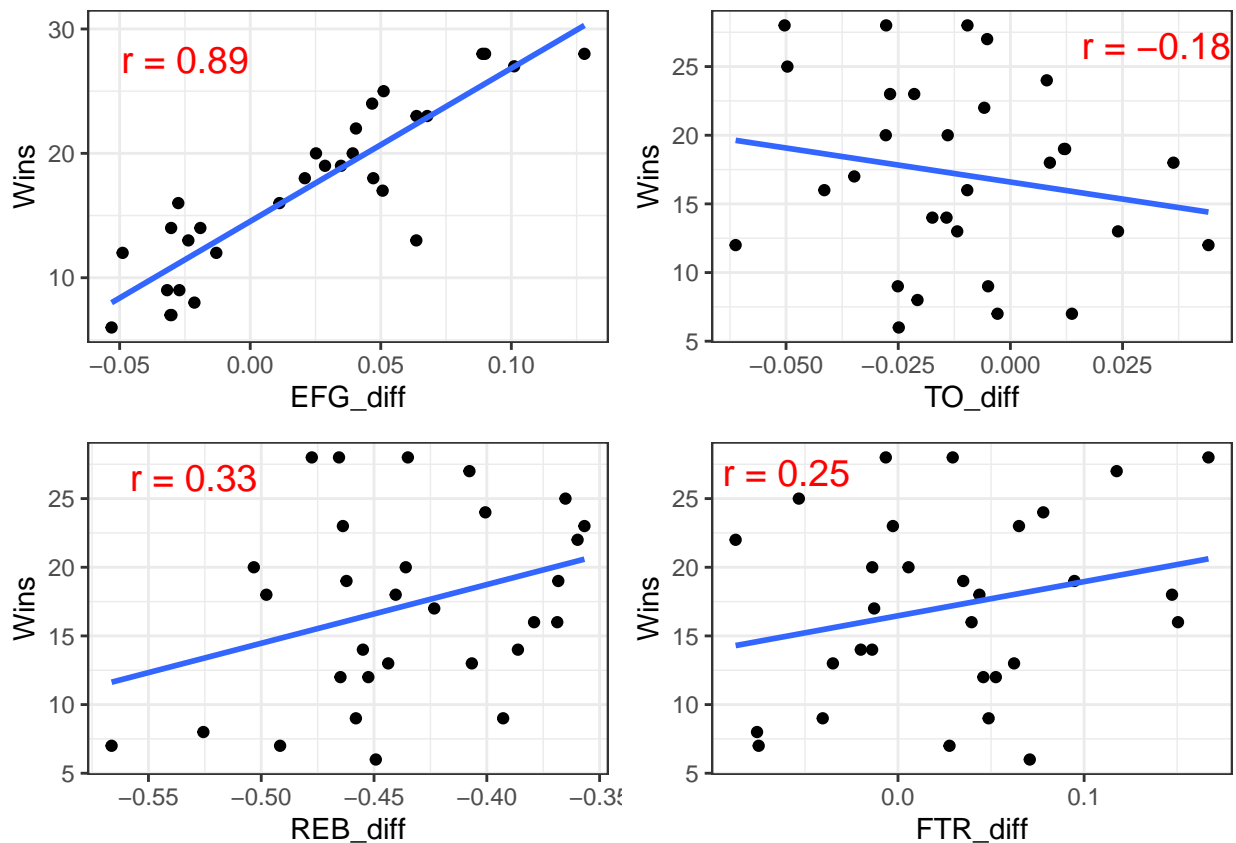
# Correlation matrix
kable(cor(deviation_tab))
```

	Wins	EFG_diff	TO_diff	REB_diff	FTR_diff
Wins	1.0000000	0.8896538	-0.1802731	0.3303412	0.2452503
EFG_diff	0.8896538	1.0000000	-0.0036477	0.1934178	0.2093006
TO_diff	-0.1802731	-0.0036477	1.0000000	-0.1185305	0.2586154
REB_diff	0.3303412	0.1934178	-0.1185305	1.0000000	0.0618440
FTR_diff	0.2452503	0.2093006	0.2586154	0.0618440	1.0000000

Here, note that our independent variables are not highly-correlated with each other. The highest correlation we see is between TO\_diff and FTR\_diff at 0.26. I postulate this is because teams that shoot a lot of free throws are “bigger” teams that draw more fouls, i.e. their lineups are taller than average, meaning they are worse ball-handlers and therefore commit more turnovers.

Statistically, the fact that these variables are relatively uncorrelated is a good thing, since the standard errors of our coefficient estimates will consequently be low. This means that our coefficient estimates will be reliable, allowing us to infer which variables most influence wins with confidence. The problem that results when two explanatory variables are highly correlated is called “collinearity,” and here is a simple discussion of the problem that anyone can understand.

Before we run our regression, let's visualize each variable's relationship with wins via scatterplots. I include the correlation of each variable with wins, denoted by “r,” and include regression lines.



The Wins vs. EFG\_diff plot... textbook. Note the high correlation there as well. Furthermore, we see a negative correlation between wins and TO\_diff—since an increase in TO\_diff means committing more turnovers on offense, this should lead to a decrease in wins. So that makes sense. We see positive correlations in the REB\_diff and FTR\_diff plots, but they are much weaker than that of EFG\_diff.

Great! We are now ready to regress wins on these four variables.

```
# Run regression on all four factors
four_factors_fit = lm(Wins ~ EFG_diff + TO_diff + REB_diff + FTR_diff, data = deviation_tab)
summary(four_factors_fit)
```

```
##
## Call:
## lm(formula = Wins ~ EFG_diff + TO_diff + REB_diff + FTR_diff,
##     data = deviation_tab)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.138 -1.355  0.649  1.858  3.718
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    21.606     4.585   4.712 7.86e-05 ***
## EFG_diff       116.124    11.034  10.525 1.13e-10 ***
## TO_diff        -52.019    22.023  -2.362  0.0263 *
## REB_diff        17.959    10.186   1.763  0.0901 .
## FTR_diff        11.150     8.247   1.352  0.1885
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.795 on 25 degrees of freedom
## Multiple R-squared:  0.8534, Adjusted R-squared:  0.83
## F-statistic: 36.39 on 4 and 25 DF,  p-value: 4.387e-10
```

There is a lot to notice here. By looking at R-squared, we see that 85% of the variability in wins can be explained by our four explanatory variables. That's pretty high. Next, we can learn a lot from the coefficients themselves. EFG\_diff is highly significant (at the 0.1% level), while TO\_diff is significant at the 5% level and REB\_diff at the 10% level. FTR\_diff is not significant. The coefficient estimates bear this out as well: in order of decreasing magnitude, we have EFG\_diff (116.1), TO\_diff (-52.0), REB\_diff (18.0), and FTR\_diff (11.1). These results imply that EFG\_diff is by far the most important factor in a college basketball team winning games.

Now we'll run a few more regressions. First, we remove FTR\_diff since it was the one insignificant variable.

```
# Remove FTR_diff
summary(lm(Wins ~ EFG_diff + TO_diff + REB_diff, data = deviation_tab))
```

```
##
## Call:
## lm(formula = Wins ~ EFG_diff + TO_diff + REB_diff, data = deviation_tab)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.1943 -1.3321  0.5922  1.6742  4.0112
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    22.288      4.629   4.815 5.48e-05 ***
## EFG_diff       119.163     10.973  10.860 3.70e-11 ***
## TO_diff        -43.981     21.540  -2.042  0.0514 .
## REB_diff        18.737     10.330   1.814  0.0813 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.839 on 26 degrees of freedom
## Multiple R-squared:  0.8427, Adjusted R-squared:  0.8246
## F-statistic: 46.44 on 3 and 26 DF,  p-value: 1.393e-10
```

Here, note that even when we exclude FTR\_diff, 84% of the variability in wins is still explained by the other three factors as opposed to 86% we saw previously. This confirms our suspicion that FTR\_diff is the least important factor in a team's success. Next, we see that TO\_diff is now only significant at the 10% level as opposed to the 5% level. I believe that this is due to the correlation of 0.26 between FTR\_diff and TO\_diff, the highest between the four explanatory variables we saw earlier. This high correlation likely caused the coefficient for TO\_diff to be a bit higher in magnitude than it deserved in the first regression, and removing FTR\_diff caused TO\_diff's coefficient to stabilize.

The last four regressions we analyze will be single variable, i.e. we regress wins on each of the four factors individually.

```
# Individual regressions. First EFG_diff.
summary(lm(Wins ~ EFG_diff, data = deviation_tab))
```

```
##
## Call:
## lm(formula = Wins ~ EFG_diff, data = deviation_tab)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.362  -2.139   0.409   2.391   4.862
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  14.5296     0.6294   23.09 < 2e-16 ***
## EFG_diff     123.1203    11.9426   10.31 4.89e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.15 on 28 degrees of freedom
## Multiple R-squared:  0.7915, Adjusted R-squared:  0.784
## F-statistic: 106.3 on 1 and 28 DF,  p-value: 4.892e-11
```

EFG\_diff alone explains 80% of variance in a team's wins and is highly significant.

```
# Second TO_diff.
summary(lm(Wins ~ TO_diff, data = deviation_tab))
```

```
##
## Call:
## lm(formula = Wins ~ TO_diff, data = deviation_tab)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11.8245  -3.9974   0.3884   5.1095  10.9344
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   16.591     1.373   12.08 1.27e-12 ***
## TO_diff      -49.565     51.108   -0.97   0.34
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.785 on 28 degrees of freedom
## Multiple R-squared:  0.0325, Adjusted R-squared: -0.002055
## F-statistic: 0.9405 on 1 and 28 DF,  p-value: 0.3404
```

TO\_diff alone explains only 3% of variation in wins and is insignificant.

```
# Third REB_diff.
summary(lm(Wins ~ REB_diff, data = deviation_tab))
```

```
##
```



```
## Call:
## lm(formula = Wins ~ REB_diff, data = deviation_tab)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.6355  -4.5956  -0.9175   4.4379  12.5768
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    35.83      10.15   3.531  0.00146 **
## REB_diff       42.74      23.08   1.852  0.07460 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.51 on 28 degrees of freedom
## Multiple R-squared:  0.1091, Adjusted R-squared:  0.07731
## F-statistic:  3.43 on 1 and 28 DF,  p-value: 0.0746
```

REB\_diff alone explains only 9% of variation in wins and is insignificant.

```
# Last FTR_diff.
summary(lm(Wins ~ FTR_diff, data = deviation_tab))
```

```
##
## Call:
## lm(formula = Wins ~ FTR_diff, data = deviation_tab)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.2271  -5.4593  -0.6365   5.4240  11.6966
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    16.466      1.328  12.397 6.88e-13 ***
## FTR_diff       24.866      18.576   1.339   0.191
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.687 on 28 degrees of freedom
## Multiple R-squared:  0.06015, Adjusted R-squared:  0.02658
## F-statistic: 1.792 on 1 and 28 DF,  p-value: 0.1915
```

As expected, FTR\_diff is insignificant and explains only 5% of variance in wins.

## Conclusion

To summarize, our regression of wins on all four factors shows we can predict wins with the formula:  $\text{predicted wins} = 21.606 + 116.124(\text{EFG\_diff}) - 52.019(\text{TO\_diff}) + 17.959(\text{REB\_diff}) + 11.150(\text{FTR\_diff})$ . However, if we had to select one of the four factors to predict or infer team success, it would undoubtedly have to be EFG\_diff. Explicitly, good shooting and/or good defense against shooters is most valuable.

Given these results, let's quantify how college basketball teams can improve. By looking at the coefficient estimates, a 1% improvement in EFG\_diff is worth 1.16 wins. This could happen through a team improving

their own EFG by 1%, improving their defense by 1% (decreasing their opponent's EFG by 1%), or a mixture of the two. Next, a 1% increase in  $TO\_diff$  is worth -0.52 wins. This makes sense, since increasing  $TO\_diff$  (turnovers committed per possession - defensive turnovers caused per possession by opponent) is obviously bad. Hence, committing 1 more turnover per 100 possessions or forcing 1 fewer turnover per 100 opponent possessions would lead to about 0.5 fewer wins. Next, a 1% increase in  $REB\_diff$  ( $ORP - DRP$ ) would lead to 0.18 more wins. This could mean grabbing 1 more offensive rebound per 100 missed shots or grabbing 1 more defensive rebound per 100 missed shots by opponent. Finally, a 1% increase in  $FTR\_diff$  would lead to 0.11 more wins, which could occur through gaining 1 free throw attempt per 100 field goal attempts or conceding 1 less free throw attempt per 100 field goal attempts by opponent.

To compare to the NBA, recall that Winston found  $EFG\_diff$  explains 71% of variance in NBA team wins,  $TO\_diff$  15%,  $REB\_diff$  6%, and  $FTR\_diff$  essentially 0%. We found corresponding percentages of 80%, 3%, 9%, and 5%. This implies that  $FTR\_diff$  in particular is more important in college than in the NBA, while  $TO\_diff$  is more important in the NBA than in the college game. Although  $EFG\_diff$  in our analysis is 9% higher than Winston's number, I'm hesitant to say it's more important in college than in the NBA because shooting in the NBA has recently become more important than it was at the time of Winston's findings.)

To devise our own weights for college basketball, we can simply divide the absolute value of each coefficient by the sum of the absolute values of all four. Doing so, we find:

1. Shooting differential (59%)
2. Turnover differential (26%)
3. Rebounding differential (9%)
4. Free throw rate differential (6%)

Overall, shooting offense and defense is undoubtedly the king of the four factors in both college and professional basketball. Meanwhile, shooting more free throws and preventing an opponent from doing so is more important in college ball, and causing more turnovers while protecting the ball better is a bit more important in the pro game. Rebounding takes on a similar level of importance at both levels of play. To speculate why we see these differences, I believe free throws are more valuable in college for a couple reasons: 1) the games are shorter, making every point more valuable, and 2) although I don't have the data to support this claim, the variance between two college teams' free throw percentages is greater than what we see between two NBA teams. As for turnovers, I presume possessions in the NBA are more valuable than those in college since 1) college players are simply worse at protecting the ball than professional players, and 2) the playing field is much more level in the NBA, meaning games are closer.

## Further Work

The main weakness of this analysis is that it does not adjust for opponent strength. For example, ACC competition is undoubtedly fiercer than WCC competition, meaning a team like Gonzaga may have unfairly influential numbers in certain areas. For instance, the fact that Gonzaga was a particularly good rebounding team in 19-20 (see the ranking table above) could cause the importance of rebounding to be inflated. A more complete work would adjust each team's numbers according to opponent strength.