

NBA

24 July 2020

Explaining Relationships Using Regression Analysis

A regression model can be used for two different purposes. It can be used to explain how one variable (y) depends on another variable (x). Indeed, here we are often trying to identify *causal* inferences, to explain how a system works. It can also be used to make *forecasts* about future outcomes.

During this week we are going to focus on the first purpose, thinking about potential causal links in the performance of professional sports teams. We are going to focus on the performance of teams across an entire season, in terms of either win percentage or league position.

The main input to any team is, of course, the players themselves. Teams compete to hire the best players and player agents seek to find the best financial deal for their clients. It is reasonable to expect, therefore, that team expenditure on player salaries should be an important factor in determining team performance. To be clear, the logic of this is NOT that paying higher salaries will make players perform better. At the top level of professional sports, all players are highly motivated, and salary probably does not play a significant motivational role. Rather, competition for players means that salaries are likely to reflect relative abilities. Better players command higher salaries, and as a result the aggregate pay of players on the team is likely to be a good predictor of team performance.

There are a number of sources for player salary data. In the North American major leagues, salary negotiations are framed by collective bargaining agreements with player unions, which often publish individual player salary data. In European soccer leagues, aggregate salary data is to be found in audited financial statements of professional clubs, which are often available to the public (notably in England). Cricket players in the Indian Premier League have their salaries determined in a public auction.

This week we are going to examine the wage-performance relationship in four different leagues - the NBA, the English Premier League, Major League Baseball and the National Hockey League. While our focus is on the role of salaries, we will also consider other factors that might be relevant, which will help us to think about some of the issues that arise in regression analysis.

We start with the NBA.

```
# As usual, we begin by loading the packages we will need

library("readxl",quietly = TRUE)
library("tidyverse",quietly = TRUE)
```

```
# Now we load the data
```

```
NBA = read_excel("NBA pay and performance.xlsx")
```

```
NBA %>% summary()
```

```
##      Team          season      wpc      salaries
## Length:210      Min.   :2012      Min.   :0.1060      Min.   : 28938902
## Class :character 1st Qu.:2013      1st Qu.:0.3780      1st Qu.: 59376595
## Mode  :character Median :2015      Median :0.5120      Median : 73728186
##              Mean  :2015      Mean   :0.4978      Mean   : 78253392
##              3rd Qu.:2017      3rd Qu.:0.6060      3rd Qu.: 95642974
##              Max.   :2018      Max.   :0.8900      Max.   :142560102
```

```
NBA %>% str()
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':   210 obs. of  4 variables:
## $ Team      : chr  "Atlanta Hawks" "Atlanta Hawks" "Atlanta Hawks" "Atlanta Hawks" ...
## $ season    : num  2012 2013 2014 2015 2016 ...
## $ wpc       : num  0.606 0.537 0.463 0.732 0.585 0.524 0.293 0.591 0.506 0.305 ...
## $ salaries: num  55503683 60437642 58841508 62487671 82337675 ...
```

If salaries reflect player ability, then the the success of a team should depend on how much more or less it pays than its rivals. However, if we look at salaries paid out in different seasons, there is clearly inflation in player salaries. We can see this if we use the `group_by()` command to look at total salaries over the seven seasons:

```
Sumsal <- NBA %>%
  group_by(season)%>%
  dplyr::summarise(salaries = sum(salaries))%>%rename(allsal = salaries)
Sumsal
```

```
## # A tibble: 7 x 2
##   season    allsal
##   <dbl>    <dbl>
## 1  2012 1438650614
## 2  2013 1837810750
## 3  2014 1976549213
## 4  2015 2153667904
## 5  2016 2523411764
## 6  2017 3159860863
## 7  2018 3343261288
```

We can see that salaries in 2018 were more than double the level in 2012, and increased consistently from year to year. This does not imply that the players were getting better from season to season. Rather, this is a reflection of the growing revenues of the NBA, and the ability of players to bargain for a more or less constant share of this growing revenue.

So, if we now want to account for team performance in terms of salaries, we need to make sure we compare like with like. What \$1 million would buy in 2012 was not the same as what it would buy in 2018. It's easy to adjust for this. We simply divide the salary of each team in each season by the total spending of all teams in that season, so that we have a measure of salary spending relative to the competition.

To do this we first use `left_join()` to add the aggregate salaries for each season to our original dataframe:

```
NBA <- left_join(NBA, Sumsal, by="season")
head(NBA)
```

```
## # A tibble: 6 x 5
```

	Team	season	wpc	salaries	allsal
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
## 1	Atlanta Hawks	2012	0.606	55503683	1438650614
## 2	Atlanta Hawks	2013	0.537	60437642	1837810750
## 3	Atlanta Hawks	2014	0.463	58841508	1976549213
## 4	Atlanta Hawks	2015	0.732	62487671	2153667904
## 5	Atlanta Hawks	2016	0.585	82337675	2523411764
## 6	Atlanta Hawks	2017	0.524	105882053	3159860863

```
tail(NBA)
```

```
## # A tibble: 6 x 5
```

	Team	season	wpc	salaries	allsal
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
## 1	Washington Wizards	2013	0.354	58296629	1837810750
## 2	Washington Wizards	2014	0.537	66615345	1976549213
## 3	Washington Wizards	2015	0.561	74590100	2153667904
## 4	Washington Wizards	2016	0.5	83380073	2523411764
## 5	Washington Wizards	2017	0.598	105042249	3159860863
## 6	Washington Wizards	2018	0.524	126304371	3343261288

We can now create a variable which we call 'relsal', which measures the share of team's salary spend in the total spending of all teams in that season:

```
# Here we create the variable 'relsal' for the NBA
```

```
NBA[, 'relsal'] = NBA[, 'salaries'] / NBA[, 'allsal']
head(NBA)
```

```
## # A tibble: 6 x 6
```

	Team	season	wpc	salaries	allsal	relsal
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
## 1	Atlanta Hawks	2012	0.606	55503683	1438650614	0.0386
## 2	Atlanta Hawks	2013	0.537	60437642	1837810750	0.0329
## 3	Atlanta Hawks	2014	0.463	58841508	1976549213	0.0298

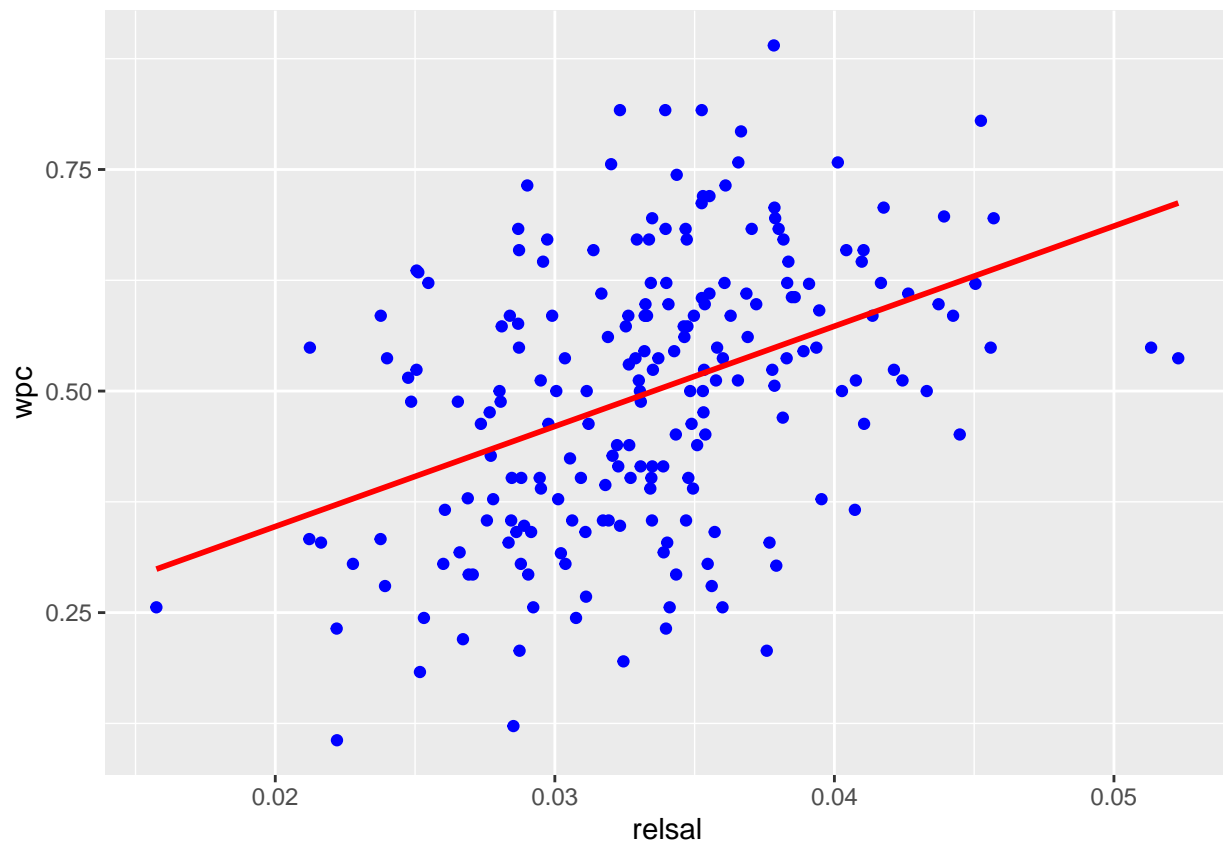
```
## 4 Atlanta Hawks    2015 0.732  62487671 2153667904 0.0290
## 5 Atlanta Hawks    2016 0.585  82337675 2523411764 0.0326
## 6 Atlanta Hawks    2017 0.524 105882053 3159860863 0.0335
```

```
tail(NBA)
```

```
## # A tibble: 6 x 6
##   Team          season  wpc  salaries    allsal relsal
##   <chr>         <dbl> <dbl>    <dbl>    <dbl> <dbl>
## 1 Washington Wizards 2013 0.354  58296629 1837810750 0.0317
## 2 Washington Wizards 2014 0.537  66615345 1976549213 0.0337
## 3 Washington Wizards 2015 0.561  74590100 2153667904 0.0346
## 4 Washington Wizards 2016 0.5    83380073 2523411764 0.0330
## 5 Washington Wizards 2017 0.598 105042249 3159860863 0.0332
## 6 Washington Wizards 2018 0.524 126304371 3343261288 0.0378
```

Before running a regression, we use `ggplot()` to look at the relationship between salaries and win percentage on a chart.

```
ggplot(data = NBA,aes(x = relsal,y = wpc )) + geom_point(color='blue') +
  geom_smooth(method = "lm", se = FALSE,color = "red")
```



It's clear from the data that there is a positive correlation between relsal and wpc, as shown by the regression line which `regplot` adds to the scatter diagram. We now run a regression

using `lm()` in order to derive the coefficients of the regression and other diagnostic statistics.

```
wpcsal1_lm = lm(formula = 'wpc ~ relsal', data = NBA)
wpcsal1_lm %>% summary()
```

```
##
## Call:
## lm(formula = "wpc ~ relsal", data = NBA)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.33885 -0.09912 -0.00710  0.09068  0.34130
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.12114     0.05809   2.086  0.0382 *
## relsal      11.30094     1.71823   6.577 3.81e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1405 on 208 degrees of freedom
## Multiple R-squared:  0.1722, Adjusted R-squared:  0.1682
## F-statistic: 43.26 on 1 and 208 DF,  p-value: 3.811e-10
```

The first things to look at in any regression are the coefficients of the explanatory variables and their statistical significance. We can see here that the coefficient on `relsal` is 11.3009. This means that every one percentage (0.01) increase in the share of the team in total salaries leads to and 11.3 x .01 increase in win percentage - that is roughly .11 , or eleven percentage points. That is a very large increase, but not that the share of salaries, as can be seen on the x axis of the chart above, ranges from around .02 (2%) to .05 (5%). Thus going from the lowest salary to the highest share (from 2% to 5%) will produce a 3 x .11 = .33 increase in win percentage - from around roughly 33% to 66%.

This estimate is statistically significant. The coefficient estimate is more than six times larger than its standard error (this ratio is called the t- statistic (6.577). The p-value ($P > |t|$) tells us the probability of observing such estimate if the true value were actually zero. The p-value here is shown as “0.000” - however, it can never be exactly zero. It is just that, in this case, it is so small that it does not register up to three decimal places. The usual standard for statistical significance is a p-value below 0.05. Clearly, in this case the estimate clearly beats that standard.

How much of the variation in team performance is captured by this `relsal`? We can see this from the R-squared coefficient which is 0.172, or 17.2%. Clearly, there is much else to team performance than salaries alone.

Is the coefficient estimate plausible? Our regression estimate is the best estimate we have of the effect, but *only under the assumption that our regression includes all of the relevant*

variables. If there are other variables which influence performance other than salary share, then our regression estimate will be biased, either upward or downward. This is the problem known as “omitted variable bias” (OVB). There is no way to be certain that OVB is not a problem, it requires good judgment and careful thought to decide on whether there are other variables that should be included.

The fact that the R-squared value was only 0.172 might give one to think that there are other factors to include, but it is always possible that the remainder is just random - the effect of luck, which no doubt plays a role in every game.

But in this case, it is possible to think of other factors that might be relevant and therefore should be included. One such is “lagged dependent variable”, which here means the value of win percentage in the previous season. While the salary level should capture many aspects of team quality, salaries are not renegotiated every year, and many aspects of team quality would have been in place in the previous season. So we can add this lagged dependent variable to our regression, and then see if this changes our estimate of the impact of salaries.

We create this variable in two stages. First, we sort the data by team, and then by season, using `arrange()` function.

Self test

Based on this model, what would be the win percentage of a team for whom the value of `realsal` was 4%?

```
NBA <- NBA %>% arrange(Team,season)
head(NBA)
```

```
## # A tibble: 6 x 6
##   Team          season  wpc  salaries    allsal realsal
##   <chr>         <dbl> <dbl>      <dbl>      <dbl>  <dbl>
## 1 Atlanta Hawks  2012 0.606  55503683 1438650614 0.0386
## 2 Atlanta Hawks  2013 0.537  60437642 1837810750 0.0329
## 3 Atlanta Hawks  2014 0.463  58841508 1976549213 0.0298
## 4 Atlanta Hawks  2015 0.732  62487671 2153667904 0.0290
## 5 Atlanta Hawks  2016 0.585  82337675 2523411764 0.0326
## 6 Atlanta Hawks  2017 0.524 105882053 3159860863 0.0335
```

```
tail(NBA)
```

```
## # A tibble: 6 x 6
##   Team          season  wpc  salaries    allsal realsal
##   <chr>         <dbl> <dbl>      <dbl>      <dbl>  <dbl>
## 1 Washington Wizards  2013 0.354  58296629 1837810750 0.0317
## 2 Washington Wizards  2014 0.537  66615345 1976549213 0.0337
## 3 Washington Wizards  2015 0.561  74590100 2153667904 0.0346
## 4 Washington Wizards  2016 0.5    83380073 2523411764 0.0330
## 5 Washington Wizards  2017 0.598 105042249 3159860863 0.0332
```

```
## 6 Washington Wizards    2018 0.524 126304371 3343261288 0.0378
```

```
NBA <- NBA %>%  
  group_by(Team)%>%  
  mutate(wpc_lag = dplyr::lag(wpc))%>%  
  ungroup()  
head(NBA)
```

```
## # A tibble: 6 x 7
```

	Team	season	wpc	salaries	allsal	realsal	wpc_lag
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
## 1	Atlanta Hawks	2012	0.606	55503683	1438650614	0.0386	NA
## 2	Atlanta Hawks	2013	0.537	60437642	1837810750	0.0329	0.606
## 3	Atlanta Hawks	2014	0.463	58841508	1976549213	0.0298	0.537
## 4	Atlanta Hawks	2015	0.732	62487671	2153667904	0.0290	0.463
## 5	Atlanta Hawks	2016	0.585	82337675	2523411764	0.0326	0.732
## 6	Atlanta Hawks	2017	0.524	105882053	3159860863	0.0335	0.585

```
tail(NBA)
```

```
## # A tibble: 6 x 7
```

	Team	season	wpc	salaries	allsal	realsal	wpc_lag
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
## 1	Washington Wizards	2013	0.354	58296629	1837810750	0.0317	0.303
## 2	Washington Wizards	2014	0.537	66615345	1976549213	0.0337	0.354
## 3	Washington Wizards	2015	0.561	74590100	2153667904	0.0346	0.537
## 4	Washington Wizards	2016	0.5	83380073	2523411764	0.0330	0.561
## 5	Washington Wizards	2017	0.598	105042249	3159860863	0.0332	0.5
## 6	Washington Wizards	2018	0.524	126304371	3343261288	0.0378	0.598

We now run our regression again, but adding wpc_lag into the regression equation:

```
wpcsal2_lm = lm(formula = 'wpc ~ wpc_lag + realsal', data = NBA)  
wpcsal2_lm %>% summary()
```

```
##  
## Call:  
## lm(formula = "wpc ~ wpc_lag + realsal", data = NBA)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -0.33064 -0.07005  0.00118  0.08533  0.26500   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  0.12925     0.05444   2.374   0.0187 *      
## wpc_lag      0.60310     0.06609   9.125  <2e-16 ***   
## realsal      2.01650     1.86111   1.083   0.2801
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1186 on 176 degrees of freedom
## (31 observations deleted due to missingness)
## Multiple R-squared:  0.4164, Adjusted R-squared:  0.4097
## F-statistic: 62.78 on 2 and 176 DF,  p-value: < 2.2e-16
```

The result of this change is quite dramatic. Not only has the coefficient on relsal fallen to only 2.0165, but it is not statistically insignificant (p_value is 0.28, greater than the critical value of .05). Last year's win percentage is highly significant and the R-squared of the regression has risen to 0.416.

It is quite usual for the lagged dependent variable to be significant in situations such as these-history matters. But should we now abandon our theory that wages influence performance? One should also be cautious about a theory that relies only on the lagged dependent variable - while history matters, one should also expect there to be specific, quantifiable factors that drive history.

While the omission of the lagged dependent variable in the first regression suggests that there may have been OVB that biased the estimate upwards, there is also the possibility that OVB can bias the estimate downwards.

Clearly, not all teams are identical, while our regression estimates thus far treat each team as if it were, in the sense that spending would affect the performance of each team in the same way, and that last year's win percentage would impact this year's in the same way. In our regression specification we want to find balance between treating each team as if it were identical, and treating each team as if it were completely unique. The truth is likely to be that there are common factors affecting all teams, but that there are also idiosyncrasies. This is often described as heterogeneity.

One way we can introduce heterogeneity is through fixed effects. Fixed effects are dummy variables. For each team there is a fixed effect, equal to one if the row relates to the team in question, and zero otherwise. Each team can have its own fixed effects. Estimation of fixed effects allows us to identify differences between the teams that are independent of the impact of salaries or of the lagged dependent variable.

Adding fixed effects is very easy in R. The variable 'Team' identifies the team names, and if we add "factor(Team)" to the regression formula, R will estimate a fixed effect for each team.

We now run the regression with the lagged dependent variable and fixed effects, to see what impact this has on the estimate of the salary coefficient.

Self test

Based on this model, what would be the win percentage of a team that (a) had 0.5 win percentage last season and (b) had a value of relsal equal to 3%?


```
wpcsal3_lm <- lm(wpc ~ wpc_lag + relsal + factor(Team),
                 data = NBA)
```

```
wpcsal3_lm %>% summary()
```

```
##
## Call:
## lm(formula = wpc ~ wpc_lag + relsal + factor(Team), data = NBA)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-0.294627	-0.067238	0.002725	0.063521	0.315694

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.2333450	0.0743186	3.140	0.00204
wpc_lag	0.2353017	0.0854354	2.754	0.00663
relsal	4.9387798	2.1020799	2.349	0.02013
factor(Team)Boston Celtics	0.0178757	0.0634989	0.282	0.77871
factor(Team)Brooklyn Nets	-0.1430698	0.0694556	-2.060	0.04117
factor(Team)Charlotte Bobcats	-0.0254028	0.0647996	-0.392	0.69561
factor(Team)Chicago Bulls	-0.0293438	0.0636228	-0.461	0.64533
factor(Team)Cleveland Cavaliers	0.0063372	0.0659487	0.096	0.92358
factor(Team)Dallas Mavericks	-0.0349739	0.0635606	-0.550	0.58299
factor(Team)Denver Nuggets	-0.0070167	0.0634787	-0.111	0.91214
factor(Team)Detroit Pistons	-0.0607367	0.0649100	-0.936	0.35096
factor(Team)Golden State Warriors	0.1682771	0.0638705	2.635	0.00932
factor(Team)Houston Rockets	0.1155091	0.0632127	1.827	0.06968
factor(Team)Indiana Pacers	-0.0168743	0.0635887	-0.265	0.79110
factor(Team)Los Angeles Clippers	0.0641937	0.0644623	0.996	0.32097
factor(Team)Los Angeles Lakers	-0.1530379	0.0673842	-2.271	0.02459
factor(Team)Memphis Grizzlies	-0.0114011	0.0640422	-0.178	0.85895
factor(Team)Miami Heat	0.0144919	0.0653730	0.222	0.82487
factor(Team)Milwaukee Bucks	-0.0502629	0.0643936	-0.781	0.43632
factor(Team)Minnesota Timberwolves	-0.0782858	0.0656635	-1.192	0.23509
factor(Team)New Orleans Hornets	-0.0585614	0.0665339	-0.880	0.38020
factor(Team)New York Knicks	-0.1144296	0.0666855	-1.716	0.08828
factor(Team)Oklahoma City Thunder	0.0649010	0.0642279	1.010	0.31393
factor(Team)Orlando Magic	-0.1411094	0.0652435	-2.163	0.03217
factor(Team)Philadelphia 76ers	-0.1193238	0.0663928	-1.797	0.07435
factor(Team)Phoenix Suns	-0.1058937	0.0644483	-1.643	0.10250
factor(Team)Portland Trail Blazers	0.0360544	0.0634828	0.568	0.57094
factor(Team)Sacramento Kings	-0.1125413	0.0657910	-1.711	0.08927
factor(Team)San Antonio Spurs	0.1328336	0.0644003	2.063	0.04091

```

## factor(Team)Toronto Raptors      0.0778177  0.0635339   1.225  0.22260
## factor(Team)Utah Jazz             0.0034318  0.0634856   0.054  0.95696
## factor(Team)Washington Wizards    -0.0009162  0.0644252  -0.014  0.98867
##
## (Intercept)                      **
## wpc_lag                          **
## relsal                           *
## factor(Team)Boston Celtics
## factor(Team)Brooklyn Nets         *
## factor(Team)Charlotte Bobcats
## factor(Team)Chicago Bulls
## factor(Team)Cleveland Cavaliers
## factor(Team)Dallas Mavericks
## factor(Team)Denver Nuggets
## factor(Team)Detroit Pistons
## factor(Team)Golden State Warriors **
## factor(Team)Houston Rockets       .
## factor(Team)Indiana Pacers
## factor(Team)Los Angeles Clippers
## factor(Team)Los Angeles Lakers    *
## factor(Team)Memphis Grizzlies
## factor(Team)Miami Heat
## factor(Team)Milwaukee Bucks
## factor(Team)Minnesota Timberwolves
## factor(Team)New Orleans Hornets
## factor(Team)New York Knicks       .
## factor(Team)Oklahoma City Thunder
## factor(Team)Orlando Magic         *
## factor(Team)Philadelphia 76ers    .
## factor(Team)Phoenix Suns
## factor(Team)Portland Trail Blazers
## factor(Team)Sacramento Kings      .
## factor(Team)San Antonio Spurs     *
## factor(Team)Toronto Raptors
## factor(Team)Utah Jazz
## factor(Team)Washington Wizards
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1094 on 147 degrees of freedom
## (31 observations deleted due to missingness)
## Multiple R-squared:  0.5851, Adjusted R-squared:  0.4976
## F-statistic: 6.686 on 31 and 147 DF,  p-value: 5.878e-16

```

There are 30 fixed effects listed in the output, one for each team. A positive fixed effect means

that in some way the team was able to perform above average, and a negative fixed effect implies below average performance. However, most of the fixed effects are not statistically significant. The significant fixed effects are: The Brooklyn Nets (negative), the Golden State Warriors (positive), The LA Lakers (negative), Orlando Magic (negative), San Antonio Spurs (positive).

Looking at our main variable of interest, *realsal*, it is clear that this variable is once again statistically significant with the addition of the fixed effects. Thus, we might conclude that the absence of the fixed effects biased the coefficient estimate downwards. The coefficient is statistically significant at the 5% level, and the value of 4.9388 implies that increasing salary share in the NBA by 1 percentage point (e.g. from 2% to 3%) will lead to an increase in win percentage of almost 5%. This is smaller than our original estimate, but also because of the presence of the lagged dependent variable, spending which increases win percentage this season will also have an effect, albeit a smaller one, on the following season. Indeed, an increase in spending today will create a ripple effect which will be discernible in performance for a number of years into the future.

Also note that the size of the lagged dependent variable is smaller once we add the fixed effects. Finally, with this third specification the R-squared of the regression has now risen to 0.585 (close to 60%), which accounts for significant fraction of the overall variation.

We should never expect to explain 100% of the variation of outcomes in sport - if we could do that then each game would be perfectly predictable - and then what would be the point of watching?

Self test

Run the regression of win percentage on *realsal* with fixed effects but without the lagged dependent variable. Compare your output results. Compare this to the previous three regressions. Which do you think is the best representation of the data. Why?

Conclusion

In this notebook we have explored the possibility of using regression analysis to explore the validity of a causal explanation of team success. That causal explanation was itself not derived from the data, but based on a theory that player quality will be reflected in salaries and therefore salaries will predict team success.

You might be wondering about why this works at all with the NBA, since the league operates a salary cap system which is intended to equalize resources among the teams. If each team spent the same amount of money on players, our theory predicts that each team can expect to win 50% of its games, and team performances will vary randomly around this mean. However, the NBA cap is a “soft cap”, meaning that there are many exemptions, so teams spend varying amounts in reality. Some leagues, such as the NFL, operate a hard cap, which strictly prohibit spending above the cap. The NFL also has a salary floor, which prevents teams from

spending a lot less than average. When looking at NFL data, therefore, it is much harder to identify the effect of wage spending on performance.

We next turn to look at the salary performance relationship in the English Premier League.