

Report

Data Wrangling and Exploratory Data Analysis

Setting up the connection

```
library(tidyverse)
library(RSQLite)
db <- DBI::dbConnect(RSQLite::SQLite(), "lahman2016.sqlite")
```

Problem 1 Using SQL, write a query to compute the total payroll and winning percentage (number of wins / number of games * 100) for each team (that is, for each teamID and yearID combination). You should include other columns that will help when performing EDA later on (e.g., franchise ids, number of wins, number of games).

```
select Salaries.yearID as year,
       Teams.teamID as team,
       (cast(sum(Teams.W) as float) / sum(Teams.G))*100 as win_percentage,
       franchName,
       sum(Salaries.salary) as payroll
from Salaries
join Teams
on Salaries.teamID = Teams.teamID
and Salaries.yearID = Teams.yearID
join TeamsFranchises
on Teams.franchID = TeamsFranchises.franchID
where year between 1990 and 2014
group by year, team
```

```
head(winpay)
```

##	year	team	win_percentage	franchName	payroll
## 1	1990	ATL	40.12346	Atlanta Braves	14555501
## 2	1990	BAL	47.20497	Baltimore Orioles	9680084
## 3	1990	BOS	54.32099	Boston Red Sox	20558333
## 4	1990	CAL	49.38272	Los Angeles Angels of Anaheim	21720000
## 5	1990	CHA	58.02469	Chicago White Sox	9491500
## 6	1990	CHN	47.53086	Chicago Cubs	13624000

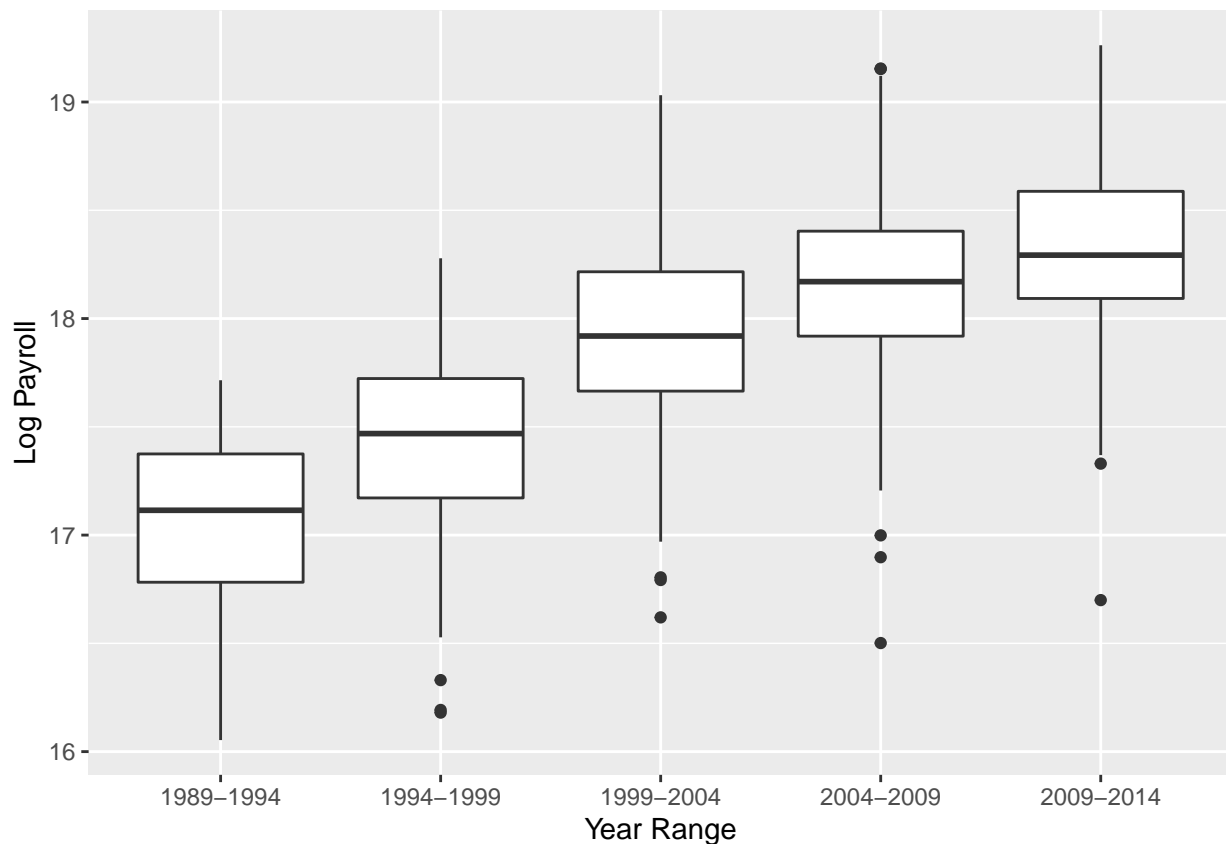
Exploratory data analysis

Payroll distribution

Problem 2. Write code to produce a plot (or plots) that shows the distribution of payrolls across teams conditioned on year (from 1990-2014). Note: you may create a single plot as long as the distributions for each year are clearly distinguishable (e.g., a single plot overlaying histograms is not OK).

```
library(ggplot2)

winpay %>%
  group_by(team) %>%
  mutate(log_payroll = log(payroll)) %>%
  mutate(yearRange = ifelse(1989 < year & year <= 1994, paste("1989-1994"),
    ifelse(1994 < year & year <= 1999, paste("1994-1999"),
      ifelse(1999 < year & year <= 2004, paste("1999-2004"),
        ifelse(2004 < year & year <= 2009, paste("2004-2009"),
          paste("2009-2014")))))) %>%
  ggplot(aes(x=yearRange, y=log_payroll)) +
  geom_boxplot() +
  labs(x="Year Range", y="Log Payroll")
```



Question 1. What statements can you make about the distribution of payrolls conditioned on time based on these plots? Remember you can make statements in terms of central tendency,

spread, etc.

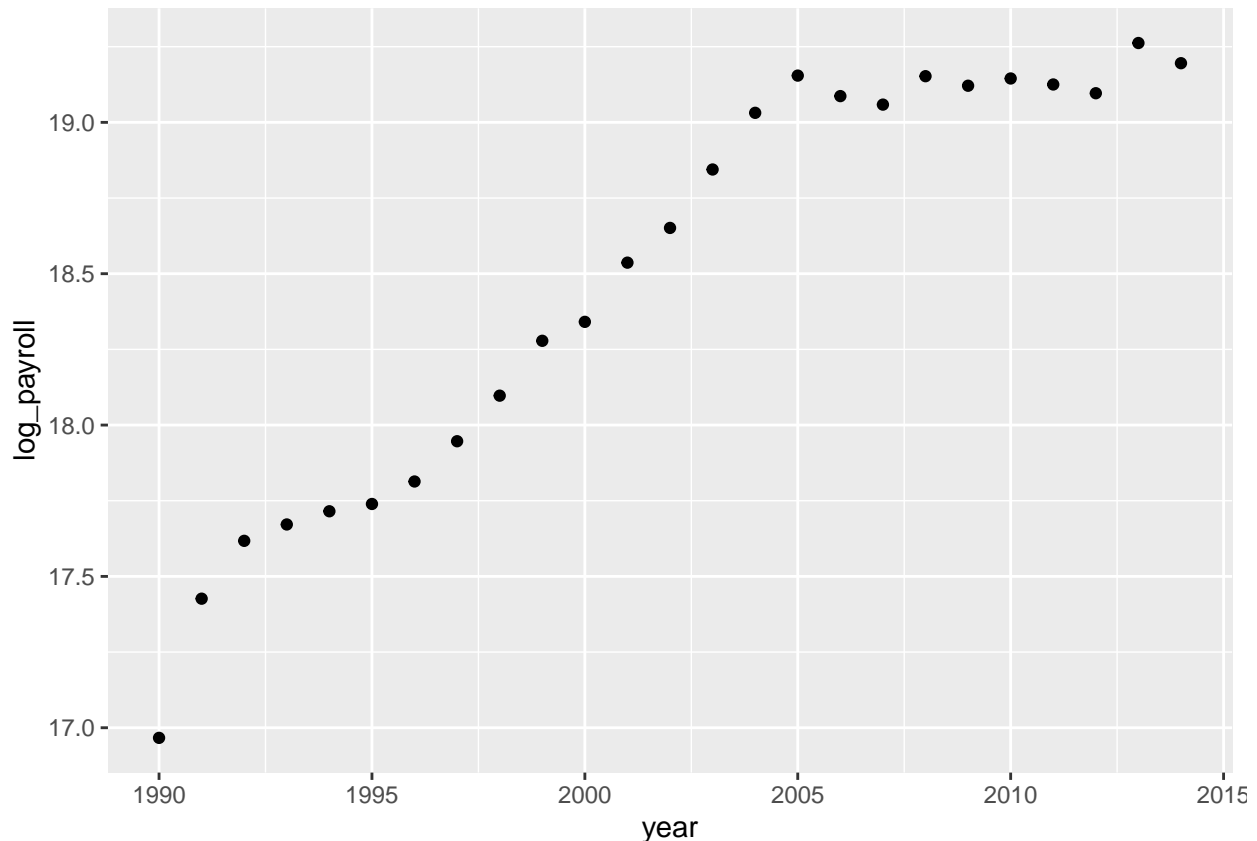
There is a clear increasing central tendency, although it seems the spread of the data has also slightly increased over time as well. It can also be seen that the upper end of the payroll scale seemed to flatten out after 2004. Most of the outliers reside on the lower end of the payroll scale, whereas at the upper end there is only one outlier.

Problem 3. Write code to produce a plot (or plots) that specifically shows at least one of the statements you made in Question 1. For example, if you make a statement that there is a trend for payrolls to decrease over time, make a plot of a statistic for central tendency (e.g., mean payroll) vs. time to show that specifically.

The statement I am going to prove is that the maximum payrolls stopped increasing significantly after 2004. I did this by graphing strictly the max payrolls between the years 1990 and 2014. The graph clearly shows a steep decline in the slope of the graph after 2005, showing that the top payrolls stopped increasing at that point.

```
select year, max(payload) as max_payload
from
  (select Salaries.yearID as year,
    Teams.teamID as team,
    sum(Salaries.salary) as payroll
  from Salaries
  join Teams
  on Salaries.teamID = Teams.teamID
  and Salaries.yearID = Teams.yearID
  where year between 1990 and 2014
  group by year, team)
group by year
```

```
max_df %>%
  mutate(log_payroll = log(max_payload)) %>%
  ggplot(aes(x=year, y=log_payroll)) +
  geom_point()
```



Correlation between payroll and winning percentage

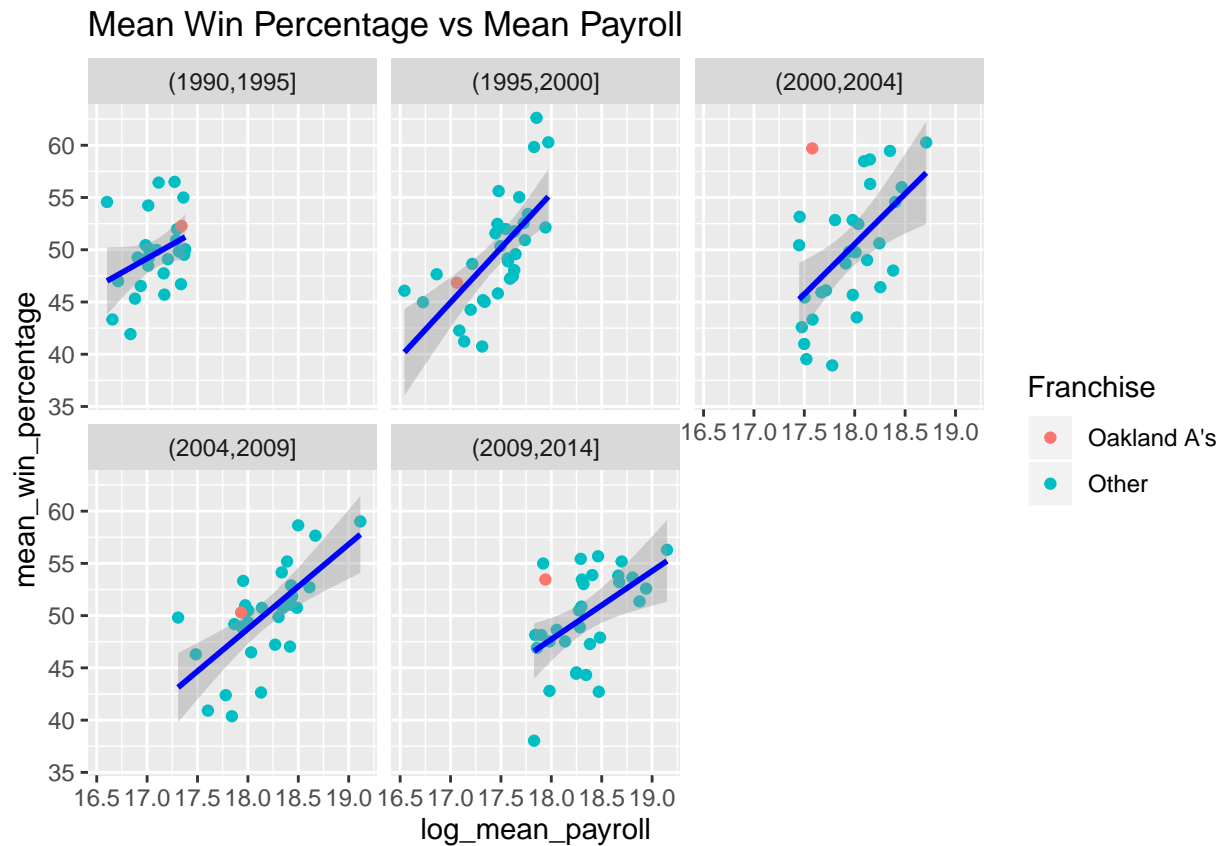
Problem 4. Write code to discretize year into five time periods (e.g., using the `cut` function with parameter `breaks=5` (in R, `bins=5` in python)) and then make a scatterplot showing mean winning percentage (y-axis) vs. mean payroll (x-axis) for each of the five time periods. You could add a regression line (using `geom_smooth(method=lm)`) in each scatter plot to ease interpretation. Note: look at the discussion on faceting in the visualization EDA lecture notes.

```
winpay$year_range <- cut(winpay$year, breaks=5)
winpay$Franchise[winpay$franchName != 'Oakland Athletics'] <- 'Other'
winpay$Franchise[winpay$franchName == 'Oakland Athletics'] <- 'Oakland'

means <- aggregate(list(winpay$payroll, winpay$win_percentage),
                    list(winpay$team, winpay$year_range), mean)
colnames(means) <- c('team', 'year_range', 'mean_payroll', 'mean_win_percentage')

means %>%
  mutate(log_mean_payroll = log(mean_payroll)) %>%
  mutate(Franchise = ifelse(team == 'OAK', 'Oakland A\'s', 'Other')) %>%
  ggplot(aes(x=log_mean_payroll, y=mean_win_percentage, color=Franchise)) +
  labs(title='Mean Win Percentage vs Mean Payroll') +
  facet_wrap(~year_range) +
```

```
geom_point() +  
geom_smooth(method=lm, color='blue')
```



Question 2. What can you say about team payrolls across these periods? Are there any teams that standout as being particularly good at paying for wins across these time periods? What can you say about the Oakland A's spending efficiency across these time periods (labeling some points in the scatterplot can help interpretation).

I can say that, as a general trend, win percentages increase as a team's payroll increases, but there's so much variance in these plots that it is hard to call this a strong correlation relationship. It seems that the time period when these two variables were most strongly correlated was the 1995 to 2000 time period. Most of the teams seem to hover either right above or right below average in terms of payroll and win percentage. The Oakland A's win percentage is exceptional. Across the five time periods, as they generally end to be either above the mean, or very close to it, except for the 1990 to 1995 time range.

Data transformations

Standardization across years

Problem 5. Write code to create a new variable in your dataset that standardizes payroll conditioned on year.

```

payroll_mew = array()
payroll_sigma = array()
for(i in 1990:2014) {
  payroll_mew[i - 1990 + 1] <- mean(filter(winpay, year==i)$payroll)
  payroll_sigma[i - 1990 + 1] <- sd(filter(winpay, year==i)$payroll)
}

winpay <- winpay %>%
  mutate(Payroll_Z_Score =
    (payroll - payroll_mew[year - 1990 + 1]) / payroll_sigma[year - 1990 + 1])

```

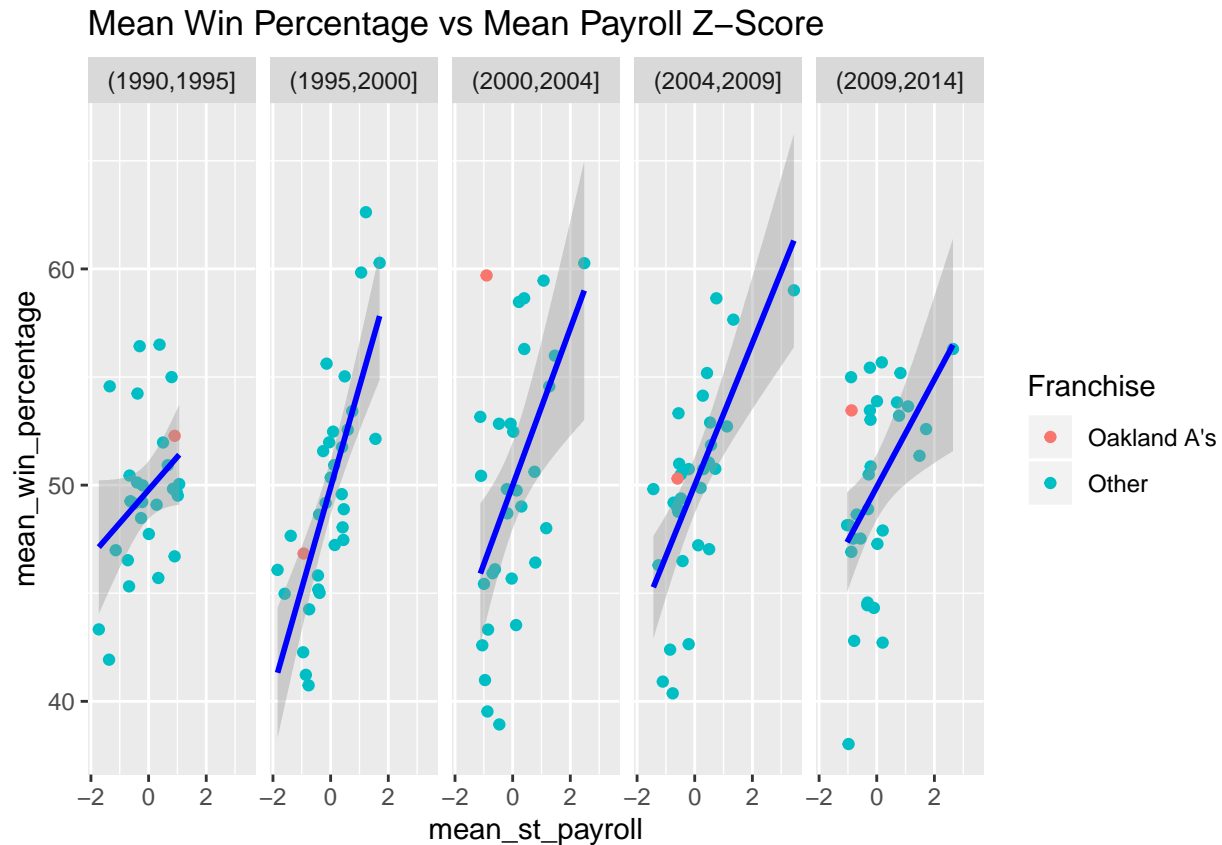
Problem 6. Repeat the same plots as Problem 4, but use this new standardized payroll variable.

```

st_means <- aggregate(list(winpay$Payroll_Z_Score, winpay$win_percentage),
  list(winpay$team, winpay$year_range), mean)
colnames(st_means) <- c('team', 'year_range', 'mean_st_payroll', 'mean_win_percentage')

st_means %>%
  mutate(Franchise = ifelse(team == 'OAK', 'Oakland A\'s', 'Other')) %>%
  ggplot(aes(x=mean_st_payroll, y=mean_win_percentage, color=Franchise)) +
  labs(title='Mean Win Percentage vs Mean Payroll Z-Score') +
  facet_grid(~year_range) +
  geom_point() +
  geom_smooth(method=lm, color='blue')

```



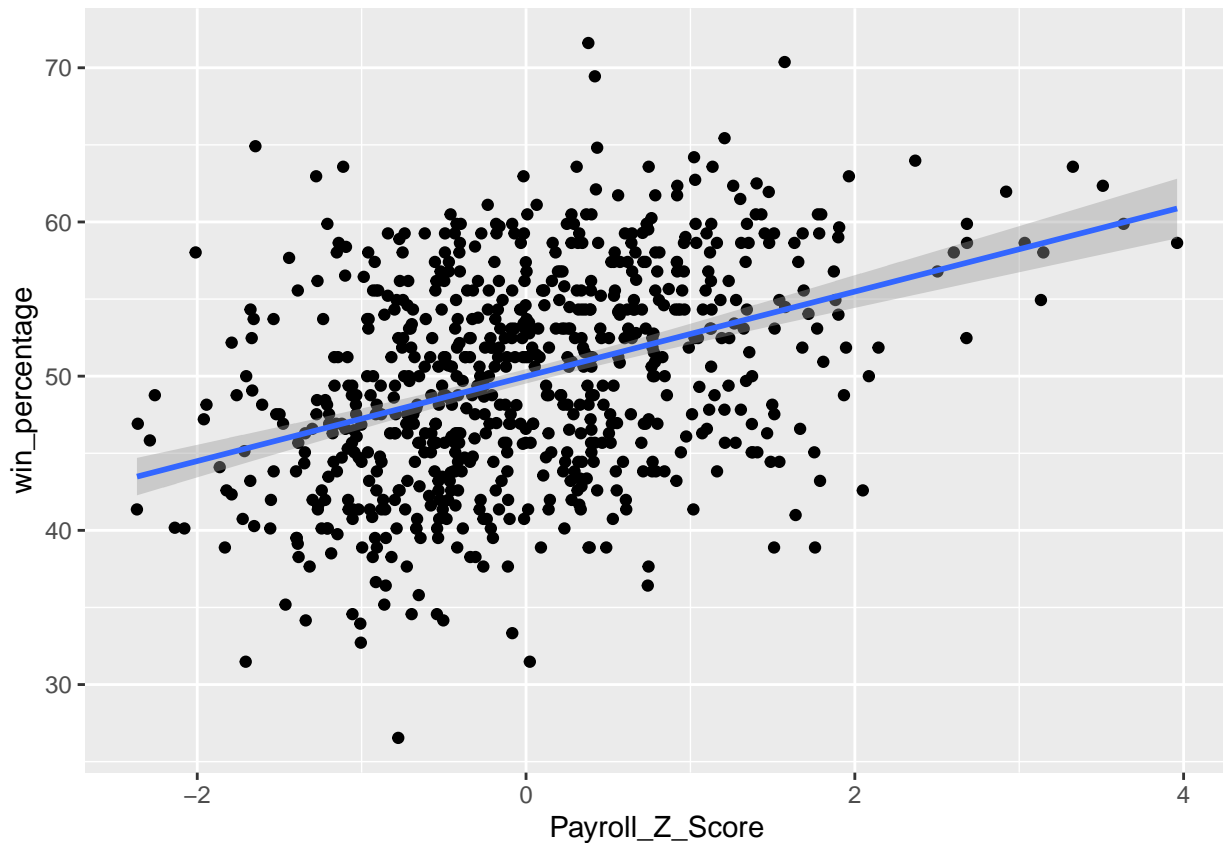
Question 3. Discuss how the plots from Problem 4 and Problem 6 reflect the transformation you did on the payroll variable. Consider data range, center and spread along with observed correlation in your discussion. Some of these change after the transformation, others don't.

The most obvious difference is the slope of the trendline has noticeably increased. The variance of the data has increased slightly, contributing to the increase in the trendline. Each graph is also centered around $x=0$, whereas the graphs in problem 4 had their own center-point making them harder to compare. Also each graph shows that at $x=0$, the win percentage is consistently 50%, implying when a team has the average payroll amount, they win 50% of the time.

Expected wins

Problem 7. Make a single scatter plot of winning percentage (y-axis) vs. standardized payroll (x-axis). Add a regression line to highlight the relationship (again using `geom_smooth(method=lm)`).

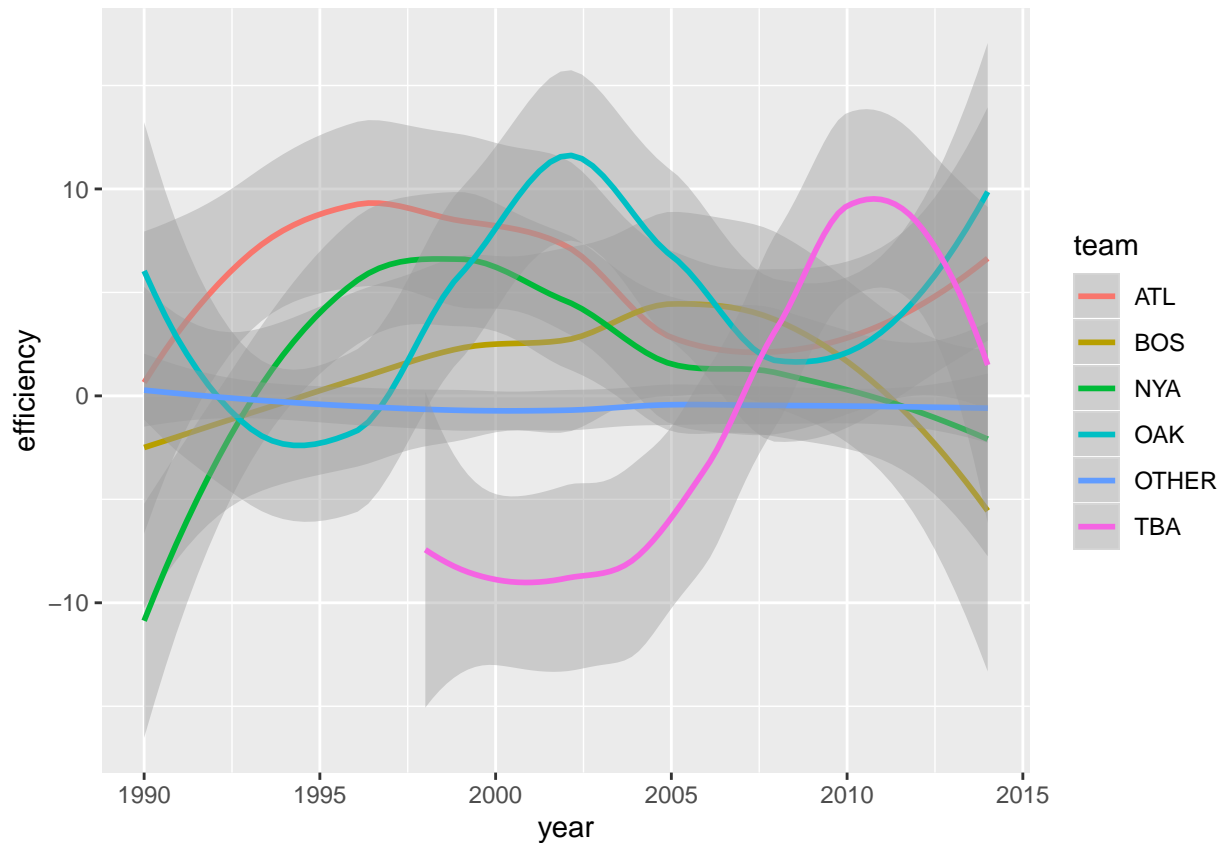
```
winpay %>%
  ggplot(aes(x=Payroll_Z_Score, y=win_percentage)) +
  geom_point() +
  geom_smooth(method=lm)
```



Spending efficiency

Problem 8. Write code to calculate spending efficiency for each team. Make a line plot with year on the x-axis and efficiency on the y-axis. A good set of teams to plot are Oakland, the New York Yankees, Boston, Atlanta and Tampa Bay (teamIDs OAK, BOS, NYA, ATL, TBA). That plot can be hard to read since there is so much year to year variation for each team. One way to improve it is to use `geom_smooth` instead of `geom_line`.

```
winpay %>%
  mutate(expected_win_pct = 50 + 2.5 * Payroll_Z_Score) %>%
  mutate(efficiency = win_percentage - expected_win_pct) %>%
  mutate(team =
    ifelse(
      team %in% c('OAK', 'BOS', 'NYA', 'ATL', 'TBA'),
      team, 'OTHER')) %>%
  ggplot(aes(x=year, y=efficiency, color=team)) +
  geom_smooth(method=loess)
```

Question 4. What can you learn from this plot compared to the set of plots you looked at in Question 2 and 3? How good was Oakland's efficiency during the Moneyball period?

This plot provides significantly more valuable information than the plots from questions 2 and 3. This plot shows which teams are worth analyzing based on their efficiency over time. For example, if I were going to start a baseball team, I would analyze Oakland's strategies for team building in 2002, where they had an efficiency of 16.3% despite consistently being at the lower end of the payroll spectrum. Oakland's efficiency during the Moneyball period was excellent. As this chart shows, Oakland started off with about average efficiency, then gradually, through the 1990s, their win percentage went up and up until it peaked in 2002. During the Moneyball period Oakland's win percentage was substantially above average and their win percentage stayed above average until roughly 2015.