# Lab 6: Principal Components Regression and Partial Least Squares

We will continue to use the `Hitters` data set in the `ISLR` package to predict `Salary` for baseball players.

```r
library(ISLR)
library(tidyverse)
library(knitr)


str(Hitters)
```

```
## 'data.frame':    322 obs. of  20 variables:
##  $ AtBat    : int  293 315 479 496 321 594 185 298 323 401 ...
##  $ Hits     : int  66 81 130 141 87 169 37 73 81 92 ...
##  $ HmRun    : int  1 7 18 20 10 4 1 0 6 17 ...
##  $ Runs     : int  30 24 66 65 39 74 23 24 26 49 ...
##  $ RBI      : int  29 38 72 78 42 51 8 24 32 66 ...
##  $ Walks    : int  14 39 76 37 30 35 21 7 8 65 ...
##  $ Years    : int  1 14 3 11 2 11 2 3 2 13 ...
##  $ CAtBat   : int  293 3449 1624 5628 396 4408 214 509 341 5206 ...
##  $ CHits    : int  66 835 457 1575 101 1133 42 108 86 1332 ...
##  $ CHmRun   : int  1 69 63 225 12 19 1 0 6 253 ...
##  $ CRuns    : int  30 321 224 828 48 501 30 41 32 784 ...
##  $ CRBI     : int  29 414 266 838 46 336 9 37 34 890 ...
##  $ CWalks   : int  14 375 263 354 33 194 24 12 8 866 ...
##  $ League   : Factor w/ 2 levels "A","N": 1 2 1 2 2 1 2 1 2 1 ...
##  $ Division : Factor w/ 2 levels "E","W": 1 2 2 1 1 2 1 2 2 1 ...
##  $ PutOuts  : int  446 632 880 200 805 282 76 121 143 0 ...
##  $ Assists  : int  33 43 82 11 40 421 127 283 290 0 ...
##  $ Errors   : int  20 10 14 3 4 25 7 9 19 0 ...
##  $ Salary   : num  NA 475 480 500 91.5 750 70 100 75 1100 ...
##  $ NewLeague: Factor w/ 2 levels "A","N": 1 2 1 2 2 1 1 1 2 1 ...
```

## 0.1 Data Processing

1. Remove records with missing values from the data (Hint: `complete.cases()` is useful)

## 0.2 Principal Components Regression

The `pcr()` function in the `pls` package can perform principal components regression.

1. Fit the PCR model using the `pcr` command. A couple tips: a) setting `scale = TRUE` will standardize your data prior to fitting the model, and b) setting `validation = TRUE` will perform 10-fold cross validation for each value of $M$.

2. Create a plot of the CV MSE (note root MSE is reported) vs. $M$.

3. When does the smallest cross-validation error occur? Which $M$ would you choose for your final model?

4. The `summary` function also provides the *percentage of variance explained* in the predictors and the response using $M$ principal components. How many principal components would we need to explain at least 80% of the variability in the predictors?

5. How much variability in $Y$ is explained for your chosen value of $M$?

## 0.3 Partial Least Squares

The `plsr()` function in the `pls` package can perform partial lest squares.

1. Fit the PLS model using the `pls` command. Again, a) setting `scale = TRUE` will standardize your data prior to fitting the model, and b) setting `validation = TRUE` will perform 10-fold cross validation for each value of $M$.

2. Create a plot of the CV MSE (note root MSE is reported) vs. $M$.

3. When does the smallest cross-validation error occur? Which $M$ would you choose for your final model?

4. How much variability in $Y$ is explained for your chosen value of $M$?

5. Discuss the two methods performed today. Which would you prefer?