

LINEAR MODEL SELECTION AND REGULARIZATION

Esam Mahdi

Data Analytics (ECMP 5005)
School of Mathematics and Statistics
Master of Engineering - Engineering Practice
Carleton University

November 10, 2023

By the end of this chapter, you should be able to do the following:

- ① Select the best subset predictors in linear regression models.
- ② Perform stepwise regression and best model selection.
- ③ Differentiate between the forward and backward stepwise techniques.
- ④ Use regularization technique: Ridge, lasso, and elastic-net regression to shrink regression coefficients and pick the most relevant features.
- ⑤ Use principle component analysis (PCA) to reduce the dimensionality of data.
- ⑥ Perform principle component regression (PCR).
- ⑦ Use R with some real-life applications.

Recall the linear regression model is given by

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon, \quad \text{where } \varepsilon \underset{\sim}{iid} \mathcal{N}(0, \sigma^2)$$

When there are too many predictor variables then we have two problems:

- 1 The prediction variance can increase a lot.
- 2 There might be a significant collinearity among some features (predictors) makes it hard to interpret the regression coefficients.

Select a subset of only relevant features will help us to reduce the prediction variance and get better interpretation and predictions.

THREE METHODS FOR FEATURE SELECTION

- *Subset selection*: Fit a regression model by least squares on a reduced subset of the p predictors that you believe to be related to the response variable.
- *Shrinkage (Regularization)*: Fit a model involving all p predictors, but some parameter estimates are shrunk towards zero by applying some penalty or constraint.
- *Dimension reduction*: Compute $M < p$ different linear combinations, or projections, of the variables and use these M projections to fit a least squares model.

Best Subset and Stepwise Model Selection

Consider the p independent variables x_1, x_2, \dots, x_p and a response y . The *best subset selection* method in a regression model :

- ① Start with the *null model*, which contains no predictors. This model has the intercept coefficient only. It simply predicts the sample mean of the response variable. Denote this model by \mathcal{M}_0
- ② For $k = 1, 2, \dots, p$:
 - ① Fit all $\binom{p}{k} = \frac{p!}{k!(p-k)!}$ models that contain exactly k predictors.
 - ② Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here *best* is defined as having the smallest RSS, or equivalently largest adjusted R^2 .
- ③ Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_k$ using cross-validated prediction error, C_p , AIC, AICc, BIC, or adjusted R^2 .

Note that this method is a computational expensive as you need to fit a large number of models, $\sum_{k=0}^p \binom{p}{k}$. For example, if $p = 10$, then you need to fit 1024 different models.

BEST SUBSET SELECTION IN R

Consider the `Hitters` data set in the **ISLR2** library. In this example, we use the `regsubsets()` function, which is a part of the package **leaps** to perform a best subset selection of relevant predictors that can be used to predict a baseball player's `Salary`. The best model is quantified using RSS.

- The syntax of the `regsubsets()` function is the same as for `lm()`.
- The `summary()` command outputs the best set of variables for each model size.
- The `summary()` returns R^2 , RSS, adjusted R^2 , C_p , and BIC (that we will explain later!).
- By default, `regsubsets()` only reports results up to the best 8 variable model.
 - The `nvmax()` option can be used in order to return as many variables as are desired.
 - The asterisk in the summary output indicates that a given variable is included in the corresponding model.
- Plotting R^2 , RSS, adjusted R^2 , C_p , and BIC for all of the models at once will help us decide which model to select.

The R code is given in the next slide, where we select the six-variable model with the lowest BIC:

$$\widehat{\text{Salary}} = 91.51 - 1.87 \text{ AtBat} + 7.60 \text{ Hits} + 3.70 \text{ Walks} + 0.64 \text{ CRBI} - 122.95 \text{ DivisionW} + 0.26 \text{ PutOuts}$$

BEST SUBSET SELECTION IN R

```
library(ISLR2); library(leaps)
head(Hitters)
anyNA(Hitters) #check if the data has any missing observations
# Count the number of missing values in Salary variable (59 "NA" observation)
sum(is.na(Hitters$Salary))
Hitters <- na.omit(Hitters) #remove the missing data
# Use nvmax() = 19 to return 19 models (default nvmax() = 8)
regfit.full <- regsubsets(Salary ~., data = Hitters, nvmax = 19)
reg.summary <- summary(regfit.full)
reg.summary
reg.summary$adjr2 # return adjusted R-squared for all models
which.max(reg.summary$adjr2) #identify which model has largest adj-R2
which.min(reg.summary$bic) #identify which model provides smallest BIC
# Plot R2, RSS, adjusted R2, Cp, and BIC for all of the models at once
par(mfrow = c(2, 2))
plot(regfit.full, scale = "r2")
plot(regfit.full, scale = "adjr2")
plot(regfit.full, scale = "Cp")
plot(regfit.full, scale = "bic")
# See the coefficient estimates associated with model 6
coef(regfit.full, 6)
```

STEPWISE REGRESSION

Stepwise Regression is the step-by-step iterative construction of a regression model that involves the selection of independent variables to be used in a final model. It involves *adding* or *removing* potential explanatory variables in succession and testing for statistical significance after each iteration (see <https://www.investopedia.com>).

Three main approaches for stepwise regression:

- ① *Forward selection.*
- ② *Backward elimination.*
- ③ *Bidirectional stepwise:* A combination of forward and backward stepwise selection, testing at each step for independent variables to be included or excluded.

FORWARD STEPWISE SELECTION

Consider the p independent variables x_1, x_2, \dots, x_p and a response y . The *forward stepwise selection* is a step-by-step iterative method for *adding* variables in a regression model:

- 1 Starts with no predictor variables in the model (*null model*)

$$y = \beta_0 + \varepsilon, \quad \text{thus} \quad \hat{y} = \hat{\beta}_0,$$

where $\hat{\beta}_0 = \bar{y}$ is the estimated point of β_0 in this case.

- 2 Add the *most significant* independent variable to this model (say x_3).

This can be achieved by implementing any of the following criteria:

- Calculate the p-value for testing the hypotheses $H_0 : \beta_j = 0$ against $H_A : \beta_j \neq 0$, for $j = 1, 2, \dots, p$ and select the *smallest p-value*.
- It provides the *highest increase* in adjusted coefficient of determination R_{adj}^2 .
- It provides the *highest drop* in *predicted residual error sum of squares (PRESS)* statistic compared to other predictors under consideration.

Now the new regression model after adding x_3 is

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_3 x_3 \quad (x_3 \text{ is added to the null model})$$

- 3 Repeat step 2 for the remaining variables until the pre-specified stopping rule (*will be explain later*) is reached or until all the independent variables are included in the model.

BACKWARD STEPWISE ELIMINATION

Consider the p independent variables x_1, x_2, \dots, x_p and a dependent y . The *backward stepwise elimination* is a step-by-step iterative method for *deleting* variables from the model:

- 1 Starts with all predictor variables in the model (*full model*)

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon, \quad \text{thus} \quad \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p,$$

where $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ are the least squares (or MLE) estimates of $\beta_0, \beta_1, \dots, \beta_p$.

- 2 Delete the *least significant* independent variable to this model (say x_3).

This can be achieved by implementing any of the following criteria:

- Calculate the p-value for testing the hypotheses $H_0 : \beta_j = 0$ against $H_A : \beta_j \neq 0$, for $j = 1, 2, \dots, p$ and select the *highest p-value* (to be deleted).
- It provides the *lowest drop* in adjusted coefficient of determination R_{adj}^2 .
- It provides the *lowest increase* in *predicted residual error sum of squares (PRESS)* statistic compared to other predictors under consideration.

Now the new regression model after deleting x_3 is

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_4 x_4 + \dots + \hat{\beta}_p x_p \quad (x_3 \text{ is removed from the model})$$

- 3 Repeat step 2 for the remaining variables until the pre-specified stopping rule (*explained in the next slides*) is reached or until no variable is left in the model.

The final fitted regression model is the model that satisfies:

- All included independent variables are significant (at pre-specified α level, say 5%) when we are testing the hypotheses $H_0 : \beta_j = 0$ versus $H_A : \beta_j \neq 0$.
- This model has the lowest *Information Criterion (IC)* (see next slide).

Thus, in stepwise regression, we stop *adding/removing* variables to our final model if *adding/removing* will:

- Provide p-values that are larger than α , which leads to not reject H_0 ;
- increase the values of IC;
- the adjusted coefficient of determination, R_{adj}^2 , value starts to decrease or not change.

THE INFORMATION CRITERION (IC)

The *Information Criterion (IC)* is a measure of goodness of fit that can help us to select the best parsimonious fitted model (i.e., the model with the least number of significant predictors).

Th popular IC are:

- *Mallow's C_p .*
- *Akaike's Information Criterion (AIC).*
- *Corrected AIC (AICc).*
- *Bayesian Information Criterion (BIC).*

The smaller the value of IC the preferred the model.

MALLOW'S C_p

$$C_p = \frac{1}{n}(\text{RSS} + 2k\hat{\sigma}^2),$$

where

- $k = p + 1$ is the number of model parameters $(\beta_0, \beta_1, \dots, \beta_p)$.
- RSS is the residual sum of squares.
- $\hat{\sigma}^2 = \text{MSE} = \text{RSS}/(n - p - 1)$ is the estimate of $\text{Var}(\varepsilon)$.

AKAIKE'S INFORMATION CRITERION (AIC)

$$\text{AIC} = 2k - 2 \log \hat{L}$$

where k denotes the number of parameters and \hat{L} is the maximized value of the *log likelihood function* for the estimated model.

Note: In the case of the linear model with Gaussian errors, maximum likelihood and least squares are the same thing, and C_p and AIC are equivalent.

CORRECTED AIC (AICc)

A bias corrected version of AIC for use when the sample size is small, or when the number of parameters estimated is a moderate to large fraction of the sample size is the *Corrected AIC (AICc)* which is defined by

$$\text{AICc} = \text{AIC} + \frac{2k^2 + 2k}{n - k - 1}$$

where n denotes the sample size and k denotes the number of parameters.

BAYESIAN INFORMATION CRITERION (BIC)

$$\text{BIC} = k \log n - 2 \log \hat{L} = \frac{1}{n} (\text{RSS} + k \hat{\sigma}^2 \log n)$$

Note:

- For model selections, there is no clear choice between AIC, C_p , and BIC.
- In general, BIC places more penalty on models with many variables. Thus, it yields to select smaller models than AIC and C_p .
- BIC is asymptotically consistent as a selection criterion. Thus, several statistician prefer to use BIC and not AIC or C_p .

- **Forward stepwise selection** is recommended to be used when the number of predictors under consideration is, p , **very large** compared with the sample size n .
- **Backward stepwise elimination** is recommended to be used when we suspect of **multicollinearity**, but can not be used when $p \geq n$.

Note: Always use the backward stepwise approach, unless the number of variables is more than the sample size, where the forward stepwise is better in this case.

The main issues with stepwise selection are:

- It does not consider all possible combination of potential predictors.
- The output of the regression coefficients, confidence intervals, p-values, and R^2 , etc. might be biased.
- It produces an unstable selection of variables, especially when we have a small sample size compared to the number of variables we want to study.
- It does not consider the causal relationship between variables.

There are many suggestions to deal with the limitations of stepwise approach:

- Split the data into training and testing sets (say K -Fold CV). Use the training data to fit the model and use the testing data to validate this model.
- Use bootstrap method to check the stability of the selection.
- Consider the shrinkage methods such as LASSO regression.

FORWARD AND BACKWARD STEPWISE SELECTION IN R

Consider again the **Hitters** data set that we studied before. We can also use the `regsubsets()` function to perform forward stepwise or backward stepwise selection, using the argument `method = "forward"` or `method = "backward"`.

```
regfit.fwd <- regsubsets(Salary ~., data = Hitters,  
  nvmax = 19, method = "forward")  
summary(regfit.fwd)  
regfit.bwd <- regsubsets(Salary ~., data = Hitters,  
  nvmax = 19, method = "backward")  
summary(regfit.bwd)
```

For instance, we see that using forward stepwise selection, the best one variable model contains only **CRBI**, and the best two-variable model additionally includes **Hits**. For this data, the best one-variable through six-variable models are each identical for best subset and forward selection. However, the best seven-variable models identified by forward stepwise selection, backward stepwise selection, and best subset selection are different. See next slide!

BEST SEVEN-VARIABLE MODELS IDENTIFIED BY FORWARD STEPWISE SELECTION, BACKWARD STEPWISE SELECTION, AND BEST SUBSET SELECTION

```

round(coef(regfit.full, 7), 3)
round(coef(regfit.fwd , 7), 3)
round(coef(regfit.bwd , 7), 3)
    
```

Coefficient	Best subset selection	Forward stepwise selection	Backward stepwise selection
Intercept	79.451	109.787	105.649
AtBat	–	–1.959	–1.976
Hits	1.283	7.450	6.757
Walks	3.227	4.913	6.056
CRBI	–	0.854	–
CAtBat	–0.375	–	–
CHits	1.496	–	–
CHmRun	1.442	–	–
CRuns	–	–	1.129
CWalks	–	–0.305	–0.716
DivisionW	–129.987	–127.122	–116.169
PutOuts	0.237	0.253	0.303

Shrinkage Methods

BIAS-VARIANCE TRADEOFF

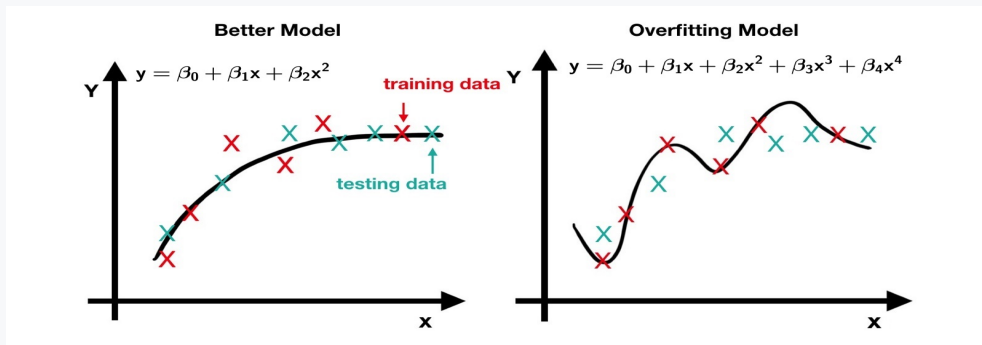


FIGURE: Over-fitting model and prediction accuracy ([Bias-variance tradeoff](#)).

Right: The model is overfit the training data. This model has poor performance with low bias (which is good) but high variance (which is bad) on the test data. Left: The model has a better performance on the test data. It decreases the variance on the cost of increasing the bias ([Bias-variance tradeoff](#)).

Recall that the multiple linear regression model can be written in the matrix notation $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ and the estimated model is

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{ols}},$$

where

$$\hat{\boldsymbol{\beta}}_{\text{ols}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}.$$

The coefficients $\hat{\boldsymbol{\beta}}_{\text{ols}}$ are obtained by minimizing the *loss function*

$$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^k x_{ij}\beta_j)^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \quad (1)$$

- Ridge, lasso, and elastic-net are used to *shrink* the least important regression coefficients (predictors) by imposing a *penalty* on their size.
- Thus, regression coefficients will be estimated by minimizing the loss function in (1) imposing some *constraint* on these coefficients.
- Adding this penalty will *reduce the variance* of parameter estimates, but it comes at the expense of some additional *bias*.
- Ridge uses *L2 regularization* technique by adding L2 penalty which is equal to the *square* of the magnitude of coefficients.
- LASSO uses *L1 regularization* technique by adding L1 penalty that is equal to the *absolute value* of the magnitude of coefficient
- Elastic-Net combines the *L1 and L2 regularization together*.
- Ridge, lasso, and elastic-net regression can be applied to *generalized linear models (glm)* including logistic regression models.
- It is best to apply ridge, lasso, and elastic net regression after *standardizing the predictors (scale)*, using the formula

$$\tilde{x}_{ij} = \frac{x_{ij}}{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2},$$

where $\bar{x}_j = 1/n \sum_{i=1}^n x_{ij}$ is the mean of the j th predictor X_j , $j = 1, 2, \dots, p$.

The estimated coefficients can be obtained by minimizing the *Lagrangian form*

$$\underbrace{\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^k x_{ij} \beta_j)^2}_{\text{Loss function}} + \underbrace{\lambda \sum_{j=1}^k \beta_j^2}_{\text{Penalty term}} = \underbrace{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}_{\text{Loss function}} + \underbrace{\lambda \boldsymbol{\beta}'\boldsymbol{\beta}}_{\text{Penalty}}, \quad (2)$$

where $0 \leq \lambda < \infty$ is the *tuning parameter*. We choose a grid of λ values and the value for which the *cross-validation* error is smallest.

The ridge regression solution of (2) is

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\mathbf{X}'\mathbf{X} + \lambda\mathbb{I})^{-1} \mathbf{X}'\mathbf{y}.$$

where \mathbb{I} is $k \times k$ identity matrix.

- The penalty term adds little bit of *bias* to the parameter estimates to *reduce the variance* of the predicted values (improve prediction accuracy).
- The intercept β_0 has been left out of the penalty term, where we estimated it by $\hat{\beta}_0 = \bar{y}$.
- When $\lambda = 0$ (no regularization), we have $\hat{\boldsymbol{\beta}}_{\text{ridge}} = \hat{\boldsymbol{\beta}}_{\text{ols}}$.
- When λ gets large asymptotically to infinity, the $\hat{\boldsymbol{\beta}}_{\text{ridge}}$ *shrinks asymptotically to zero (but never be exactly zero)*. Thus, ridge regression does not perform variable selection!

LEAST ABSOLUTE SHRINKAGE AND SELECTION OPERATOR (LASSO)

The estimated coefficients can be obtained by minimizing

$$\underbrace{\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^k x_{ij} \beta_j)^2}_{\text{Loss function}} + \underbrace{\lambda \sum_{j=1}^k |\beta_j|}_{\text{Penalty term}}, \quad (3)$$

where $0 \leq \lambda < \infty$ is the *tuning parameter* that can be chosen by *cross-validation* methods, same as in ridge regression. The lasso regression solution of (3) can be obtained using a *quadratic programming technique*.

- The penalty term adds little bit of *bias* to the parameter estimates to *reduce the variance* of the predicted values (improve prediction accuracy).
- When $\lambda = 0$ (no regularization), we have $\hat{\beta}_{\text{lasso}} = \hat{\beta}_{\text{ols}}$.
- Making λ sufficiently large can cause some of the coefficients (least important predictors) to be *exactly zero*. Thus, LASSO is a method for feature selection.
- Thus, lasso regression effectively selects the most important predictors for predicting outcome (*better than stepwise regression*).

The estimated coefficients can be obtained by minimizing

$$\underbrace{\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^k x_{ij} \beta_j)^2}_{\text{Loss function}} + \lambda \underbrace{\sum_{j=1}^k ((1 - \alpha) \beta_j^2 + \alpha |\beta_j|)}_{\text{Penalty mixed term}}, \quad (4)$$

where $0 \leq \lambda < \infty$ is the *tuning parameter* and $0 \leq \alpha \leq 1$ is the *elastic-net mixing parameter* which determines the mix of the penalties for considering both ridge and lasso.

- The penalty term adds little bit of *bias* to the parameter estimates to *reduce the variance* of the predicted values (improve prediction accuracy).
- λ and α can be chosen by *cross-validation* methods.
- When $\alpha = 0$, we have ridge regression, i.e. $\hat{\beta}_{\text{elastic-net}} = \hat{\beta}_{\text{ridge}}$.
- When $\alpha = 1$, we have lasso regression, i.e. $\hat{\beta}_{\text{elastic-net}} = \hat{\beta}_{\text{lasso}}$.
- When $\lambda = 0$ (no regularization), we have $\hat{\beta}_{\text{elastic-net}} = \hat{\beta}_{\text{ols}}$.
- λ is chosen by cross-validation.

RIDGE VS LASSO REGRESSION

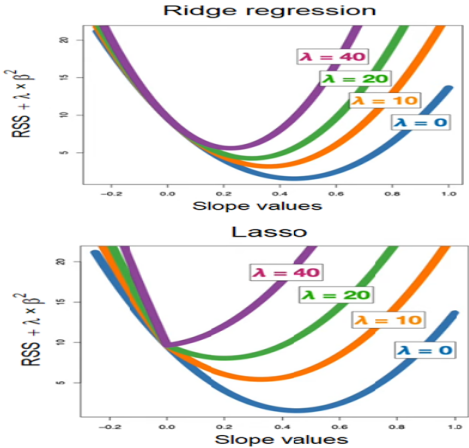
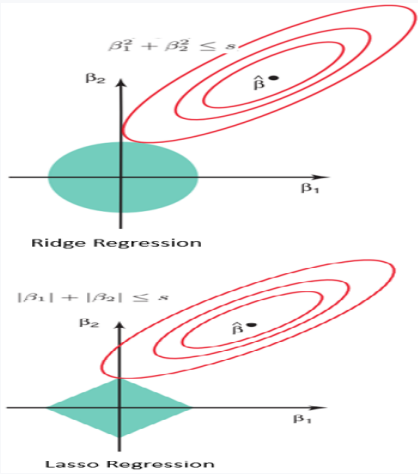


FIGURE: Lasso Regression can shrink parameter values towards 0, but Ridge Regression can not.

RIDGE, LASSO, AND ELASTIC-NET REGRESSION IN R

THE `GLMNET()` FUNCTION IN THE PACKAGE `glmnet`

- The `glmnet()` function in the package `glmnet` can be used to fit the ridge, lasso, elastic-net models, and more.
- More information about the R package `glmnet` can be found [here](#).
- The syntax of this function is slightly different from other model-fitting functions as we must pass in an x matrix as well as a y vector, and we do not use the $y \sim x$ syntax.
- The `glmnet()` function has an `alpha` argument that determines what type of model is fit. If `alpha = 0` then a ridge regression model is fit, and if `alpha = 1` then a lasso model is fit.
- By default the argument `lambda` in the `glmnet()` is used for an automatically selected range of λ values. However, you can implement the function over a grid of values.
- By default, the `glmnet()` function standardizes the variables so that they are on the same scale, which is always recommended. To turn off this default setting, use the argument `standardize = FALSE`.

See Section 6.5 Lab: Linear Models and Regularization Methods from the textbook.

THE `TRAIN()` FUNCTION IN THE PACKAGE `caret`

See the R Markdown example file uploaded in the Brightspace of this course.

Dimension Reduction Methods

PRINCIPLE COMPONENT ANALYSIS (PCA)

The idea in *principal component analysis (PCA)* is to combine p predictor variables into a smaller set of variables ($M \ll p$), which are weighted linear combinations of the original set.

- The set of M variables (denotes *principal components*) explains most of the variability of the full set of variables.
 - The M principal components (PC) are uncorrelated.
 - The first PC accounts for as much of the variability in the data as possible.
 - Each succeeding PC accounts for as much of the remaining variability as possible.

KEY TERMS FOR PRINCIPAL COMPONENTS ANALYSIS

- *Principal component*: A linear combination that transforms a large number of correlated variables into a smaller number of *uncorrelated variables* called principal components.
- *Loadings*: The weights (range from -1 and 1) for each original variable that transform the predictors into the components.
- *Screeplot*: A plot shows the number of principle components on the X –axis against the proportion of the variance explained by these components on the Y –axis.

PRINCIPAL COMPONENTS REGRESSION (PCR)

For two predictor variables, X_1 and X_2 , there are two principal components Z_1 and Z_2 that are defined as a linear combination of the original predictor variables:

$$Z_i = \phi_{1i}X_1 + \phi_{2i}X_2, \quad i = 1, 2$$

The coefficients (weights) ϕ_{1i}, ϕ_{2i} are known as the *loadings*. The PCA algorithm chooses ϕ 's to capture as much variance of the data as possible. The first principal component, Z_1 , captures the largest variance of the data. The second principal component, Z_2 , captures the remaining variation.

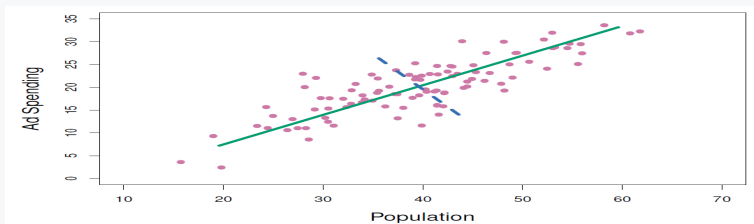


FIGURE: The **green solid line** indicates the first principal component, and the **blue dashed line** indicates the second principal component.

Now, we can fit a **simple linear regression model** using one PCA:

$$Y = \theta_0 + \theta_1 Z_1 + \varepsilon = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon,$$

where $\beta_0 = \theta_0$, $\beta_1 = \theta_1 \phi_{11}$ and $\beta_2 = \theta_1 \phi_{21}$.

PRINCIPAL COMPONENTS REGRESSION (PCR)

For a data with p predictor variables, X_1, X_2, \dots, X_p , the M principle components, Z_1, Z_2, \dots, Z_M , where $M \ll p$, are defined as linear combinations of the original p predictors as follows:

$$Z_m = \phi_{1m}X_1 + \phi_{2m}X_2 + \dots + \phi_{pm}X_p = \sum_{j=1}^p \phi_{jm}X_j, \quad m = 1, 2, \dots, M.$$

We can then use the least squares method to fit the [multiple linear regression model](#)

$$y_i = \theta_0 + \sum_{m=1}^M \theta_m z_{im} + \varepsilon_i, \quad i = 1, \dots, n \quad (5)$$

Notice that from (5),

$$\begin{aligned} \sum_{m=1}^M \theta_m z_{im} &= \sum_{m=1}^M \theta_m \sum_{j=1}^p \phi_{jm} x_{ij} = \sum_{j=1}^p \sum_{m=1}^M \theta_m \phi_{jm} x_{ij} \\ &= \sum_{j=1}^p \beta_j x_{ij}, \quad \text{where} \quad \beta_j = \sum_{m=1}^M \theta_m \phi_{jm} \end{aligned}$$

Hence, (5) can be seen as a special case of the original linear regression model and dimension reduction serves to constrain the estimated β_j coefficients.

PRINCIPAL COMPONENTS REGRESSION (PCR) IN R: HITTERS BASEBALL DATASET

Principal components regression (PCR) can be performed using the `pcr()` function, which is part of the **ppls** package.

- The syntax for the `pcr()` function is similar to that for `lm()` with some additional parameters.
- The default argument `center = TRUE` is used to perform mean centering and setting the argument `scale = TRUE` will standardize predictor variables.
- Setting the argument `validation = "CV"` tells R to calculate 10-fold cross-validation error for each possible value of the number of principal components used. Also note that you can specify `validation = "LOOCV"` instead to perform *leave-one-out cross-validation*.
- To determine the number of principal components worth keeping, you can examine the resulting fit using the syntax `summary()`.

We now apply the PCR to the `Hitters` data, in order to predict the `Salary` (response variable). See the R Markdown example file uploaded in the Brightspace of this course!

CALCULATE THE PCA USING THE VARIANCE–COVARIANCE MATRIX

Let \mathbf{X} denotes the data set of n rows and p features: $\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$.

- Calculate the *eigenvalue-eigenvector* pairs $(\lambda_1, \phi_1), (\lambda_2, \phi_2), \dots, (\lambda_p, \phi_p)$ of the variance-covariance matrix Σ (defined in Theorem 1, later), where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ and $\phi_i = (\phi_{1i}, \phi_{2i}, \dots, \phi_{pi})'$.
- The i th principal component is given by

$$Z_i = \phi_i' \mathbf{X} = \phi_{1i} X_1 + \phi_{2i} X_2 + \cdots + \phi_{pi} X_p, \quad i = 1, 2, \dots, p. \quad (6)$$

- Given a $n \times p$ data set and (7), we define the *scores* of the first principle component (PC) as

$$z_{i1} = \phi_{1i} x_{i1} + \phi_{2i} x_{i2} + \cdots + \phi_{pi} x_{ip}, \quad \text{where } \phi_{1i}, \dots, \phi_{pi} \text{ are the } \textit{loading} \text{ of the first PC.}$$

- We choose ϕ_1 (loading/weights vector) so that $\text{Var}(Z_1) = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j1} x_{ij} \right)^2 = \lambda_1$ is maximized subject to the constraint that $\phi_1' \phi_1 = \sum_{j=1}^p \phi_{j1}^2 = 1$.
- The second PC corresponds to the choice of ϕ_2 so that $\text{Var}(Z_2) = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j2} x_{ij} \right)^2 = \lambda_2$ is maximized subject to $\phi_2' \phi_2 = \sum_{j=1}^p \phi_{j2}^2 = 1$ and $\text{Cov}(Z_2, Z_1) = 0$.
- The i th PC corresponds to the choice of ϕ_i so that $\text{Var}(Z_i) = \lambda_i$ is maximized subject to $\phi_i' \phi_i = 1$ and $\text{Cov}(Z_i, Z_j) = 0$ for $j = 1, \dots, i-1$.
- One can show that the total variance of the original variables equals to the sum of the eigenvalues (total sum of variances of principle components). That is $\sum_{i=1}^p \text{Var}(X_i) = \lambda_1 + \lambda_2 + \cdots + \lambda_p = \sum_{i=1}^p \text{Var}(Z_i)$.
- The *proportion* of the total variance explained by the i th PC is $\lambda_i / \sum_{i=1}^p \lambda_i$.

CALCULATE THE PCA USING THE CORRELATION MATRIX

Usually, we *standardized the predictors* (column means and variances of \mathbf{X} are 0 and 1, respectively):

- Mathematically, we use the formula $\tilde{x}_{ji} = \frac{x_{ji} - \bar{x}_i}{\sqrt{S_{ii}}}$, where

$$S_{ii} = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_i)^2}, \quad \bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ji}, \quad i = 1, \dots, p.$$

- In R, we use the syntax: `scale(x, center = TRUE, scale = TRUE)`.

In this case, we calculate the *eigenvalue-eigenvector* pairs $(\lambda_1, \phi_1), (\lambda_2, \phi_2), \dots, (\lambda_p, \phi_p)$ of the correlation matrix ρ , where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ and $\phi_i = (\phi_{1i}, \phi_{2i}, \dots, \phi_{pi})'$ denotes the i^{th} principal component *loading vector*.

- The i th principal component is given by

$$Z_i = \phi_i' \tilde{\mathbf{X}} = \phi_{1i} \tilde{X}_1 + \phi_{2i} \tilde{X}_2 + \dots + \phi_{pi} \tilde{X}_p, \quad i = 1, 2, \dots, p. \quad (7)$$

- Total sum of variances of principle components. That is $\sum_{i=1}^p \text{Var}(\tilde{X}_i) = \sum_{i=1}^p \text{Var}(Z_i) = p$.
- The *proportion* of the total variance explained by the i^{th} PC is λ_i/p .

THEOREM 1

If $Z_1 = \phi_1' \mathbf{X}$, $Z_2 = \phi_2' \mathbf{X}$, \dots , $Z_p = \phi_p' \mathbf{X}$ are the principal components obtained from the covariance matrix

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{bmatrix} = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \cdots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \cdots & \text{Cov}(X_2, X_p) \\ \vdots & \vdots & \vdots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \cdots & \text{Var}(X_p) \end{bmatrix},$$

then

$$\rho_{Z_i, X_k} = \frac{\phi_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}, \quad i, k = 1, 2, \dots, p$$

are the *correlation coefficients* between the components Z_i and the variables X_k . Here $(\lambda_1, \phi_1), (\lambda_2, \phi_2), \dots, (\lambda_p, \phi_p)$ are the eigenvalue-eigenvector pairs for Σ , $\sigma_{kk} = \text{Var}(X_k)$ and

$$\sigma_{ik} = \text{Cov}(X_i, X_k) = \frac{1}{n} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k).$$

THEOREM 2

If $Z_1 = \phi_1' \tilde{\mathbf{X}}$, $Z_2 = \phi_2' \tilde{\mathbf{X}}$, \dots , $Z_p = \phi_p' \tilde{\mathbf{X}}$ are the principal components obtained from the correlation matrix ρ then

$$\rho_{Z_i, \tilde{X}_k} = \phi_{ik} \sqrt{\lambda_i} \quad i, k = 1, 2, \dots, p$$

are the *correlation coefficients* between the components Z_i and the variables \tilde{X}_k . Here $(\lambda_1, \phi_1), (\lambda_2, \phi_2), \dots, (\lambda_p, \phi_p)$ are the eigenvalue-eigenvector pairs for ρ .

PCA OBTAINED FROM VARIANCE-COVARIANCE MATRIX

Suppose the random variables X_1 , X_2 and X_3 have the variance-covariance matrix

$$\Sigma = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \text{Cov}(X_1, X_3) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \text{Cov}(X_2, X_3) \\ \text{Cov}(X_3, X_1) & \text{Cov}(X_3, X_2) & \text{Var}(X_3) \end{bmatrix} = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

The eigenvalue-eigenvector pairs can be obtained by solving $|\Sigma - \lambda I| = 0$, where I is identity matrix:

$$\lambda_1 = 5.83, \quad \phi'_1 = [.383, -.924, 0]$$

$$\lambda_2 = 2.00, \quad \phi'_2 = [0, 0, 1]$$

$$\lambda_3 = 0.17, \quad \phi'_3 = [.924, .383, 0]$$

Therefore, the principal components become

$$Z_1 = \phi'_1 \mathbf{X} = .383X_1 - .924X_2$$

$$Z_2 = \phi'_2 \mathbf{X} = X_3$$

$$Z_3 = \phi'_3 \mathbf{X} = .924X_1 + .383X_2$$

PCA OBTAINED FROM COVARIANCE MATRIX (CONT.)

The variable X_3 is one of the principal components, because it is uncorrelated with the other two variables. One can show that $\text{Var}(Z_i) = \lambda_i$ for $i = 1, 2, 3$ and $\text{Cov}(Z_i, Z_k) = 0$ for $i \neq k$. For example,

$$\begin{aligned}\text{Var}(Z_1) &= \text{Var}(.383X_1 - .924X_2) \\ &= (.383)^2 \text{Var}(X_1) + (-.924)^2 \text{Var}(X_2) \\ &\quad + 2(.383)(-.924) \text{Cov}(X_1, X_2) \\ &= .147(1) + .854(5) - .708(-2) \\ &= 5.83 = \lambda_1\end{aligned}$$

$$\begin{aligned}\text{Cov}(Z_1, Z_2) &= \text{Cov}(.383X_1 - .924X_2, X_3) \\ &= .383 \text{Cov}(X_1, X_3) - .924 \text{Cov}(X_2, X_3) \\ &= .383(0) - .924(0) = 0\end{aligned}$$

It is also readily apparent that

$$\sum_{i=1}^3 \text{Var}(X_i) = 1 + 5 + 2 = \lambda_1 + \lambda_2 + \lambda_3 = 5.83 + 2.00 + .17 = \sum_{i=1}^3 \text{Var}(Z_i)$$

PCA OBTAINED FROM COVARIANCE MATRIX (CONT.)

The proportion of total variance accounted for by the first principal component is $\lambda_1 / (\lambda_1 + \lambda_2 + \lambda_3) = 5.83/8 = .73$. Further, the first two components account for a proportion $(5.83 + 2)/8 = .98$ of the population variance. In this case, the components Z_1 and Z_2 could replace the original three variables with little loss of information. Next, using Theorem 1, we obtain

$$\rho_{Z_1, X_1} = \frac{\phi_{11}\sqrt{\lambda_1}}{\sqrt{\sigma_{11}}} = \frac{.383\sqrt{5.83}}{\sqrt{1}} = .925$$
$$\rho_{Z_1, X_2} = \frac{\phi_{12}\sqrt{\lambda_1}}{\sqrt{\sigma_{22}}} = \frac{-.924\sqrt{5.83}}{\sqrt{5}} = -.998$$

Notice here that the variable X_2 , with coefficient $-.924$, receives the greatest weight (load) in the component Z_1 . It also has the largest correlation (in absolute value) with Z_1 . The correlation of X_1 , with Z_1 , $.925$, is almost as large as that for X_2 , indicating that the variables are about equally important to the first principal component. The relative sizes of the coefficients of X_1 and X_2 suggest, however, that X_2 contributes more to the determination of Z_1 than does X_1 . Since, in this case, both coefficients are reasonably large and they have opposite signs, we would argue that both variables aid in the interpretation of Z_1 .

Finally,

$$\rho_{Z_2, X_1} = \rho_{Z_2, X_2} = 0 \quad \text{and} \quad \rho_{Z_2, X_3} = \frac{\sqrt{\lambda_2}}{\sqrt{\sigma_{33}}} = \frac{\sqrt{2}}{\sqrt{2}} = 1 \quad (\text{as it should})$$

The remaining correlations can be neglected, since the third component is unimportant.

PRINCIPAL COMPONENTS IN R: HITTERS BASEBALL DATASET

Consider the `Hitters` baseball dataset from the **ISLR2** package, which contains various information about 263 professional players on 20 variables. In this example, we perform a PCA on the two predictors `Hits` (number of hits in 1986) and `Runs` (number of runs in 1986).

```
df <- ISLR2::Hitters[,c("Hits", "Runs")] ## PCA using prcomp() function
## PCA using princomp() function      pca <- prcomp(df, scale = FALSE)
pca <- princomp(df, score = F)        pca
pca$loadings

##                                     ## Standard deviations (1, ..., p=2):
## [1] 52.497396  8.906177
##                                     ##
##                                     ## Rotation (n x k) = (2 x 2):
##                                     ##          PC1          PC2
## Hits  0.8812558 -0.4726396
## Runs  0.4726396  0.8812558

##                                     ##
##                                     ##          Comp.1  Comp.2
## SS loadings          1.0    1.0
## Proportion Var      0.5    0.5
## Cumulative Var      0.5    1.0

summary(pca)

## Importance of components:
##                                     ##          Comp.1    Comp.2
## Standard deviation    52.4158146  8.89233667
## Proportion of Variance 0.9720241  0.02797591
## Cumulative Proportion 0.9720241  1.00000000
```

The weights for `Hits` and `Runs` for the first principal component are 0.881 and 0.473 and for the second principal component they are 0.473 and -0.881 :

$Z_1 = 0.881 \text{ Hits} + 0.473 \text{ Runs}$, $Z_2 = 0.473 \text{ Hits} - 0.881 \text{ Runs}$
The first principal components explain the variability around 97% and its captures the majority of the variability.

PRINCIPAL COMPONENTS IN R: HITTERS BASEBALL DATASET

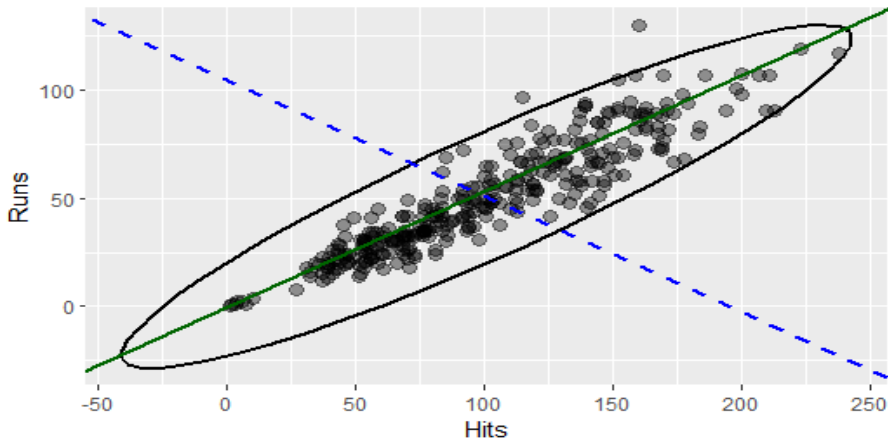


FIGURE: The green solid line indicates the first principal component, and the blue dashed line indicates the second principal component.

PCA OBTAINED FROM CORRELATION MATRIX

Consider the covariance matrix

$$\Sigma = \begin{bmatrix} 1 & 4 \\ 4 & 100 \end{bmatrix}$$

and the derived correlation matrix

$$\rho = \begin{bmatrix} \frac{1}{\sqrt{1}\sqrt{1}} & \frac{4}{\sqrt{1}\sqrt{100}} \\ \frac{4}{\sqrt{1}\sqrt{100}} & \frac{100}{\sqrt{100}\sqrt{100}} \end{bmatrix} = \begin{bmatrix} 1 & .4 \\ .4 & 1 \end{bmatrix}$$

The eigenvalue-eigenvector pairs from Σ are

$$\begin{aligned} \lambda_1 &= 100.16, & \phi_1' &= [.040, .999] \\ \lambda_2 &= .84, & \phi_2' &= [.999, -.040] \end{aligned}$$

Similarly, the eigenvalue-eigenvector pairs from ρ are

$$\begin{aligned} \lambda_1 &= 1 + \rho = 1.4, & \phi_1' &= [.707, .707] \\ \lambda_2 &= 1 - \rho = .6, & \phi_2' &= [.707, -.707] \end{aligned}$$

The respective principal components using Σ are:

$$\begin{aligned} Z_1 &= .040X_1 + .999X_2 \\ Z_2 &= .999X_1 - .040X_2 \end{aligned}$$

PCA OBTAINED FROM CORRELATION MATRIX (CONT.)

The respective principal components using ρ : are:

$$\begin{aligned}Z_1 &= .707\tilde{X}_1 + .707\tilde{X}_2 = .707\left(\frac{X_1 - \mu_1}{1}\right) + .707\left(\frac{X_2 - \mu_2}{10}\right) \\ &= .707(X_1 - \mu_1) + .0707(X_2 - \mu_2) \\ Z_2 &= .707\tilde{X}_1 - .707\tilde{X}_2 = .707\left(\frac{X_1 - \mu_1}{1}\right) - .707\left(\frac{X_2 - \mu_2}{10}\right) \\ &= .707(X_1 - \mu_1) - .0707(X_2 - \mu_2)\end{aligned}$$

Because of its large variance, X_2 completely dominates the first principal component determined from Σ . Moreover, this first principal component explains a proportion

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{100.16}{101} = .992$$

of the total population variance.

When the variables X_1 and X_2 are standardized, however, the resulting variables contribute equally to the principal components determined from ρ . Using Theorem 2, we obtain

$$\rho_{Z_1, \tilde{X}_1} = \phi_{11} \sqrt{\lambda_1} = .707\sqrt{1.4} = .837$$

and

$$\rho_{Z_2, \tilde{X}_1} = \phi_{21} \sqrt{\lambda_1} = .707\sqrt{1.4} = .837$$

PCA OBTAINED FROM CORRELATION MATRIX (CONT.)

In this case, the first principal component explains a proportion

$$\frac{\lambda_1}{p} = \frac{1.4}{2} = .7$$

of the total (standardized) population variance.

Most strikingly, we see that the relative importance of the variables to, for instance, the first principal component is greatly affected by the standardization.

When the first principal component obtained from ρ is expressed in terms of X_1 and X_2 , the relative magnitudes of the weights .707 and .0707 are in direct opposition to those of the weights .040 and .999 attached to these variables in the principal component obtained from Σ .

HOW MANY COMPONENTS TO CHOOSE?

In order to reduce the dimension of the data, you must decide how many principal components to select? Below some methods that you can use:

- Method 1 (ad hoc rule): by plotting the *screeplot*. A scree plot shows how much variation each PC captures from the data. An ideal curve should be steep, then bends at an "elbow" — this is the numbers of components you need to select — and after that flattens out.

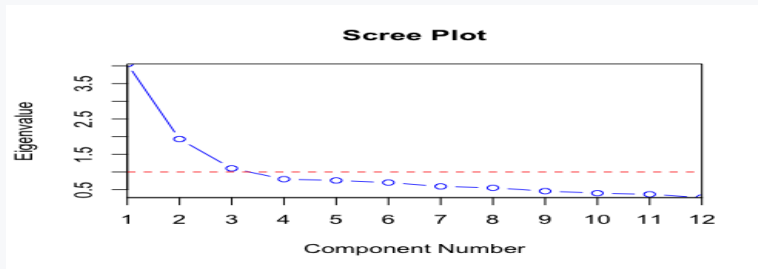


FIGURE: Source: Wikipedia. Just PC 1 – 3 might be enough to describe the data.

- Method 2 (ad hoc rule): Select the top components such that the cumulative variance exceeds a given threshold (say 85%).
- Method 3 (formal rule): Use cross-validation methods.

HITTERS BASEBALL DATASET

```
df <- as.data.frame(scale(Hitters[,c(1:13,16:18)]))  
pca <- princomp(df)  
par(mfrow=c(1,2))  
screplot(pca); screplot(pca, type = "lines")
```

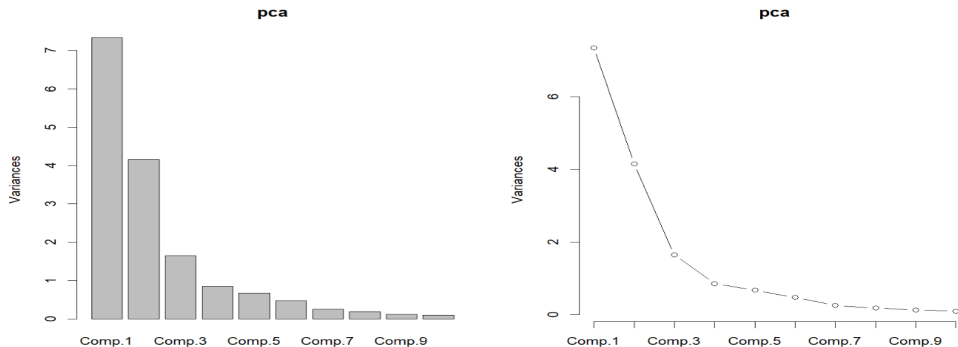


FIGURE: A screeplot for a PCA of hitters baseball data. Just PC 1 – 6 might be enough to describe the variability in the data.

You need to

- Read and understand the examples and code in Section 6.5 Lab: Linear Models and Regularization Methods from the textbook *"An Introduction to Statistical Learning: With Applications in R"*
- Solve question 5 in Exercise 6.6.
- Solve questions 8 to 11 in Exercise 6.6.