

Principal Component Regression

Nate Wells

Math 243: Stat Learning

November 16th, 2020

Outline

In today's class, we will . . .

- Discuss Principal Component Analysis as a means of dimensionality reduction for regression
- Implement PCR in R

Section 1

Principal Component Regression

Dimensionality Reduction

Suppose you collect a sample of n observations on p predictors X_1, \dots, X_p , where p is relatively large. Suppose further that some of the predictors are correlated with one another.

Dimensionality Reduction

Suppose you collect a sample of n observations on p predictors X_1, \dots, X_p , where p is relatively large. Suppose further that some of the predictors are correlated with one another.

- Any predictive model for a response Y based on all of the correlated variables will underperform due to instability in parameter estimates.

Dimensionality Reduction

Suppose you collect a sample of n observations on p predictors X_1, \dots, X_p , where p is relatively large. Suppose further that some of the predictors are correlated with one another.

- Any predictive model for a response Y based on all of the correlated variables will underperform due to instability in parameter estimates.

It may be difficult to fit complex models accurately, given limited number of observations compared to predictors.

Dimensionality Reduction

Suppose you collect a sample of n observations on p predictors X_1, \dots, X_p , where p is relatively large. Suppose further that some of the predictors are correlated with one another.

- Any predictive model for a response Y based on all of the correlated variables will underperform due to instability in parameter estimates.

It may be difficult to fit complex models accurately, given limited number of observations compared to predictors.

- If p is larger than n , it may not be possible to fit certain models to the data (for example MLR models cannot be used)

Dimensionality Reduction

Suppose you collect a sample of n observations on p predictors X_1, \dots, X_p , where p is relatively large. Suppose further that some of the predictors are correlated with one another.

- Any predictive model for a response Y based on all of the correlated variables will underperform due to instability in parameter estimates.

It may be difficult to fit complex models accurately, given limited number of observations compared to predictors.

- If p is larger than n , it may not be possible to fit certain models to the data (for example MLR models cannot be used)

One solution is to perform variable selection and drop some less useful predictors.

Dimensionality Reduction

Suppose you collect a sample of n observations on p predictors X_1, \dots, X_p , where p is relatively large. Suppose further that some of the predictors are correlated with one another.

- Any predictive model for a response Y based on all of the correlated variables will underperform due to instability in parameter estimates.

It may be difficult to fit complex models accurately, given limited number of observations compared to predictors.

- If p is larger than n , it may not be possible to fit certain models to the data (for example MLR models cannot be used)

One solution is to perform variable selection and drop some less useful predictors.

- But dropping variables completely loses possible valuable information.

Dimensionality Reduction

Suppose you collect a sample of n observations on p predictors X_1, \dots, X_p , where p is relatively large. Suppose further that some of the predictors are correlated with one another.

- Any predictive model for a response Y based on all of the correlated variables will underperform due to instability in parameter estimates.

It may be difficult to fit complex models accurately, given limited number of observations compared to predictors.

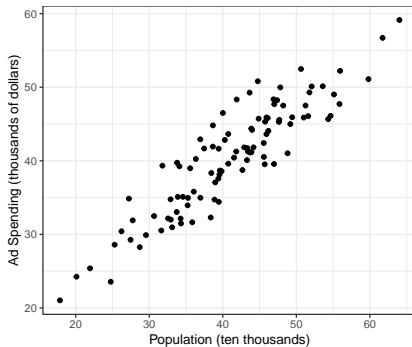
- If p is larger than n , it may not be possible to fit certain models to the data (for example MLR models cannot be used)

One solution is to perform variable selection and drop some less useful predictors.

- But dropping variables completely loses possible valuable information.
- Instead, we can combine variables into new ones that adequately describe the variance in the data, and drop those that have limited utility in explaining that variance.

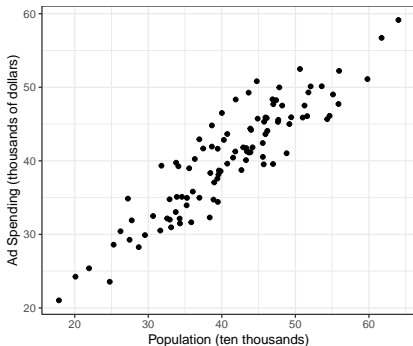
PCA Overview

Consider the relationship between campaign ad spending and population size for 100 cities:



PCA Overview

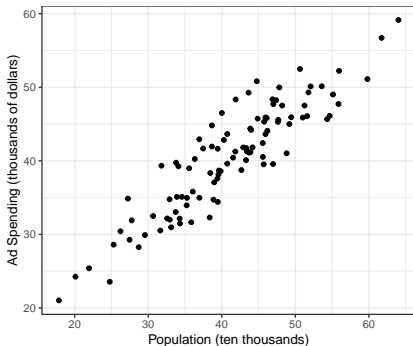
Consider the relationship between campaign ad spending and population size for 100 cities:



What are the approximate standard deviations of ad spending and population?

PCA Overview

Consider the relationship between campaign ad spending and population size for 100 cities:

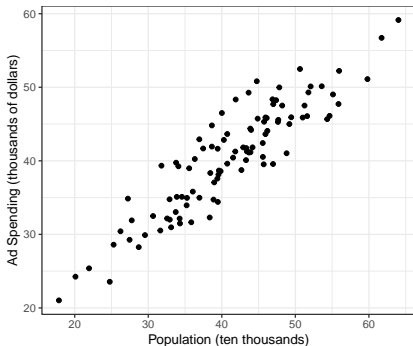


What are the approximate standard deviations of ad spending and population?

```
##      sd_Pop      sd_Ad  
## 1  8.981994  7.418227
```

PCA Overview

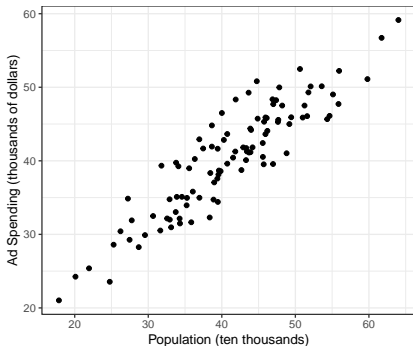
Consider the relationship between campaign ad spending and population size for 100 cities:



But how much of the variation in ad spending is just due to variation in population?

PCA Overview

Consider the relationship between campaign ad spending and population size for 100 cities:



But how much of the variation in ad spending is just due to variation in population?

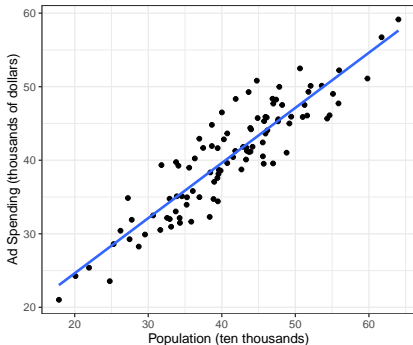
```
##           R_sq  
## 1 0.8238886
```

PCA Overview

Can we find a line along which the observations vary the most?

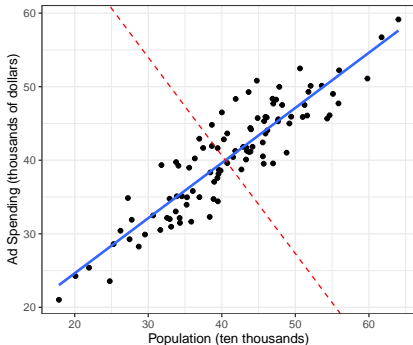
PCA Overview

Can we find a line along which the observations vary the most?



PCA Overview

How much variation occurs perpendicular to this line?



PCA Definition

The *first principal component* Z_1 is the direction along which there is the greatest variability in the data.

PCA Definition

The *first principal component* Z_1 is the direction along which there is the greatest variability in the data.

- That is, if we project the observations onto this line, the resulting projected observations would have the greatest possible variance.

PCA Definition

The *first principal component* Z_1 is the direction along which there is the greatest variability in the data.

- That is, if we project the observations onto this line, the resulting projected observations would have the greatest possible variance.
- Projecting a point onto a line amounts to finding the location on the line closest to the given point.

PCA Definition

The *first principal component* Z_1 is the direction along which there is the greatest variability in the data.

- That is, if we project the observations onto this line, the resulting projected observations would have the greatest possible variance.
- Projecting a point onto a line amounts to finding the location on the line closest to the given point.

We can express the first principal component as a linear combination of the centered predictors $X_i - \bar{X}_i$, where $\phi_{i1} \in \mathbb{R}$ and $\phi_{11}^2 + \dots + \phi_{p1}^2 = 1$:

PCA Definition

The *first principal component* Z_1 is the direction along which there is the greatest variability in the data.

- That is, if we project the observations onto this line, the resulting projected observations would have the greatest possible variance.
- Projecting a point onto a line amounts to finding the location on the line closest to the given point.

We can express the first principal component as a linear combination of the centered predictors $X_i - \bar{X}_i$, where $\phi_{i1} \in \mathbb{R}$ and $\phi_{11}^2 + \dots + \phi_{p1}^2 = 1$:

$$Z_1 = \phi_{11}(X_1 - \bar{X}_1) + \phi_{21}(X_2 - \bar{X}_2) + \dots + \phi_{p1}(X_p - \bar{X}_p)$$

PCA Definition

The *first principal component* Z_1 is the direction along which there is the greatest variability in the data.

- That is, if we project the observations onto this line, the resulting projected observations would have the greatest possible variance.
- Projecting a point onto a line amounts to finding the location on the line closest to the given point.

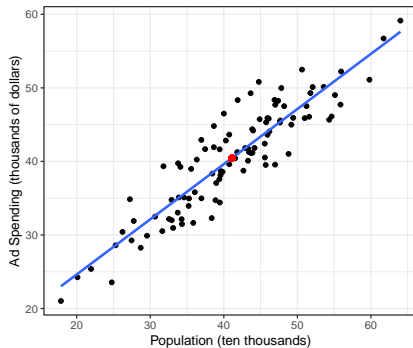
We can express the first principal component as a linear combination of the centered predictors $X_i - \bar{X}_i$, where $\phi_{i1} \in \mathbb{R}$ and $\phi_{11}^2 + \dots + \phi_{p1}^2 = 1$:

$$Z_1 = \phi_{11}(X_1 - \bar{X}_1) + \phi_{21}(X_2 - \bar{X}_2) + \dots + \phi_{p1}(X_p - \bar{X}_p)$$

- Alternatively, we could express Z_1 as an affine linear combination of the predictors themselves (affine meaning including a constant term)

PCA Example

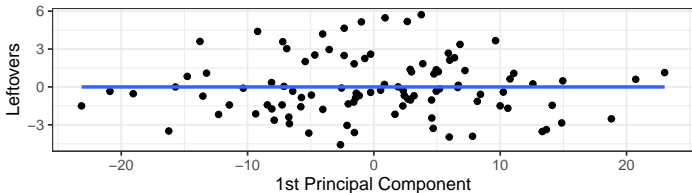
The first principal component



$$Z_1 = 0.8(\text{Pop} - 41.1) + 0.6(\text{Ad} - 40.4)$$

PCA Example

What is leftover?



Other Principal Components

In general, if we have p predictors, we can compute p distinct principal components:
 Z_1, Z_2, \dots, Z_p .

Other Principal Components

In general, if we have p predictors, we can compute p distinct principal components: Z_1, Z_2, \dots, Z_p .

The second principal component Z_2 is a linear combination of the centered variables that is

- uncorrelated with the first principal component
- has the largest variance subject to this constraint.

Other Principal Components

In general, if we have p predictors, we can compute p distinct principal components: Z_1, Z_2, \dots, Z_p .

The second principal component Z_2 is a linear combination of the centered variables that is

- uncorrelated with the first principal component
- has the largest variance subject to this constraint.

For the case when $p = 2$, the 2nd principal component corresponds to the line perpendicular to the line for the 1st principal component.

Other Principal Components

In general, if we have p predictors, we can compute p distinct principal components: Z_1, Z_2, \dots, Z_p .

The second principal component Z_2 is a linear combination of the centered variables that is

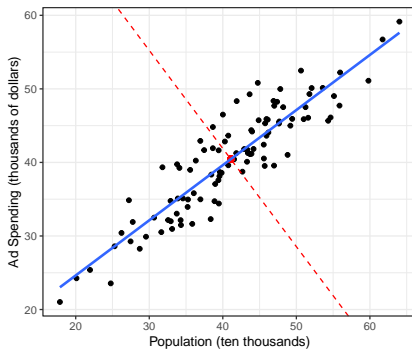
- uncorrelated with the first principal component
- has the largest variance subject to this constraint.

For the case when $p = 2$, the 2nd principal component corresponds to the line perpendicular to the line for the 1st principal component.

Generally, the k th principal component is obtained by finding a linear combination of centered variables that is uncorrelated with all previous principal components, and has the largest variance subject to this constraint.

PCA Example

The second principal component



$$Z_2 = 0.6(\text{Pop} - 41.1) - 0.8(\text{Ad} - 40.4)$$

Principal Component Regression

The PCR approach to linear regression constructs the first M principal components Z_1, \dots, Z_M of a data set with p predictors (so $M \leq p$), and then uses these as predictors in a linear regression model.

Principal Component Regression

The PCR approach to linear regression constructs the first M principal components Z_1, \dots, Z_M of a data set with p predictors (so $M \leq p$), and then uses these as predictors in a linear regression model.

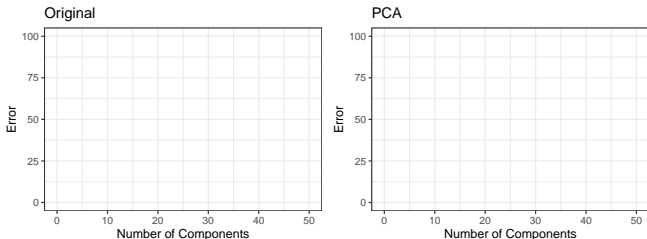
- Goal: Use a small number of predictors which explain most of the variability in the data set, as well as their relationship to the response.

Principal Component Regression

The PCR approach to linear regression constructs the first M principal components Z_1, \dots, Z_M of a data set with p predictors (so $M \leq p$), and then uses these as predictors in a linear regression model.

- Goal: Use a small number of predictors which explain most of the variability in the data set, as well as their relationship to the response.

In general, PCR tends to produce linear models with higher accuracy than models fit with the original predictors.



Principal Component Regression in R

We can use the `pcr` function in the `pls` library to quickly perform PCR in R.

Principal Component Regression in R

We can use the `pcr` function in the `pls` library to quickly perform PCR in R.

The `Hitters` data set from the `ISLR` package contains `Salary` and 18 other predictors for 263 baseball players

```
set.seed(1)
library(pls)
my_pcr <- pcr( Salary ~ ., data = Hitters, scale = T, validation = "CV")
```

Principal Component Regression in R

We can use the `pcr` function in the `pls` library to quickly perform PCR in R.

The `Hitters` data set from the `ISLR` package contains `Salary` and 18 other predictors for 263 baseball players

```
set.seed(1)
library(pls)
my_pcr <- pcr( Salary ~ ., data = Hitters, scale = T, validation = "CV")
```

- Setting `scale = T` standardizes each predictor
- Setting `validation = "CV"` causes `pcr` to compute the 10-fold CV error for each value of M (number of principal components used)

PCR Results

summary(my_pcr)

```
## Data:      X dimension: 263 19
## Y dimension: 263 1
## Fit method: svdpc
## Number of components considered: 19
##
## VALIDATION: RMSEP
## Cross-validated using 10 random segments.
##      (Intercept) 1 comps 2 comps 3 comps 4 comps 5 comps 6 comps
## CV              452    352.5   351.6   352.3   350.7   346.1   345.5
## adjCV           452    352.1   351.2   351.8   350.1   345.5   344.6
##      7 comps 8 comps 9 comps 10 comps 11 comps 12 comps 13 comps
## CV       345.4   348.5   350.4   353.2   354.5   357.5   360.3
## adjCV    344.5   347.5   349.3   351.8   353.0   355.8   358.5
##      14 comps 15 comps 16 comps 17 comps 18 comps 19 comps
## CV       352.4   354.3   345.6   346.7   346.6   349.4
## adjCV    350.2   352.3   343.6   344.5   344.3   346.9
##
## TRAINING: % variance explained
##      1 comps 2 comps 3 comps 4 comps 5 comps 6 comps 7 comps 8 comps
## X       38.31   60.16   70.84   79.03   84.29   88.63   92.26   94.96
## Salary  40.63   41.58   42.17   43.22   44.90   46.48   46.69   46.75
##      9 comps 10 comps 11 comps 12 comps 13 comps 14 comps 15 comps
## X       96.28   97.26   97.98   98.65   99.15   99.47   99.75
## Salary  46.86   47.76   47.82   47.85   48.10   50.40   50.55
##      16 comps 17 comps 18 comps 19 comps
## X       99.89   99.97   99.99   100.00
## Salary  53.01   53.85   54.61   54.61
```

PCR Results

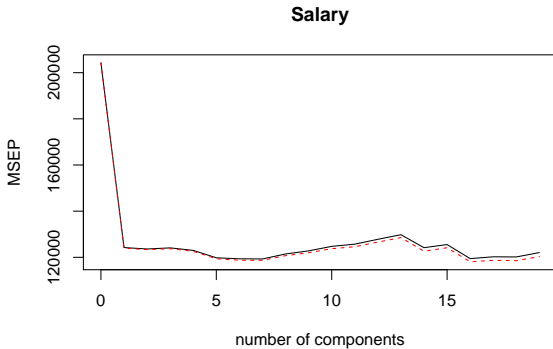
summary(my_pcr)

```
## Data:      X dimension: 263 19
## Y dimension: 263 1
## Fit method: svdpc
## Number of components considered: 19
##
## VALIDATION: RMSEP
## Cross-validated using 10 random segments.
##      (Intercept) 1 comps 2 comps 3 comps 4 comps 5 comps 6 comps
## CV              452   352.5   351.6   352.3   350.7   346.1   345.5
## adjCV           452   352.1   351.2   351.8   350.1   345.5   344.6
##      7 comps 8 comps 9 comps 10 comps 11 comps 12 comps 13 comps
## CV       345.4   348.5   350.4   353.2   354.5   357.5   360.3
## adjCV    344.5   347.5   349.3   351.8   353.0   355.8   358.5
##      14 comps 15 comps 16 comps 17 comps 18 comps 19 comps
## CV       352.4   354.3   345.6   346.7   346.6   349.4
## adjCV    350.2   352.3   343.6   344.5   344.3   346.9
##
## TRAINING: % variance explained
##      1 comps 2 comps 3 comps 4 comps 5 comps 6 comps 7 comps 8 comps
## X       38.31   60.16   70.84   79.03   84.29   88.63   92.26   94.96
## Salary  40.63   41.58   42.17   43.22   44.90   46.48   46.69   46.75
##      9 comps 10 comps 11 comps 12 comps 13 comps 14 comps 15 comps
## X       96.28   97.26   97.98   98.65   99.15   99.47   99.75
## Salary  46.86   47.76   47.82   47.85   48.10   50.40   50.55
##      16 comps 17 comps 18 comps 19 comps
## X       99.89   99.97   99.99   100.00
## Salary  53.01   53.85   54.61   54.61
```

- Note: pcr reports RSE, so values need to be squared to get MSE.

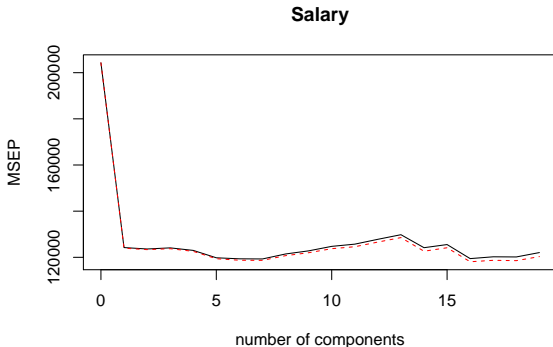
Validation Plot

```
validationplot(my_pcr, val.type = "MSEP")
```



Validation Plot

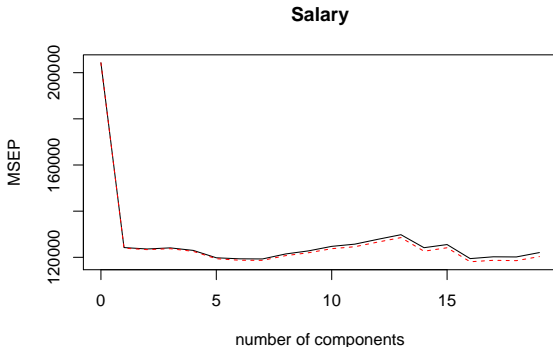
```
validationplot(my_pcr, val.type = "MSEP")
```



- Note: The smallest CV error occurs at $M = 16$ (which is close to the maximum number of predictors $p = 19$.)

Validation Plot

```
validationplot(my_pcr, val.type = "MSEP")
```



- Note: The smallest CV error occurs at $M = 16$ (which is close to the maximum number of predictors $p = 19$.)
- However, a relatively low CV error is also obtained at $M = 6$, suggesting fewer components are sufficient

Comparative Performance

Live coding. A .Rmd file will be available on course website after class