# Statistical Modeling

CH.9 - Working with Collinear Data

SS 2021 || Prof. Dr. Buchwitz

# Outline

**1** **Organizational Information**

**2** Multicollinearity

**3** Principal Components

**4** Principal Component Regression

**5** Ridge Regression

# Course Contents

| Session | Topic |
|---------|-------|
| 1 | Simple Linear Regression |
| 2 | Multiple Linear Regression |
| 3 | Regression Diagnostics |
| 4 | Qualitative Variables as Predictors |
| 5 | Transformation of Variables |
| 6 | Weighted Least Squares |
| 7 | Correlated Errors |
| 8 | Analysis of Collinear Data |
| 9 | Working with Collinear Data |
| 10 | Variable Selection Procedures |
| 11 | Logistic Regression |
| 12 | Further Topics |

# Outline

## Introduction

- When multicollinearity is present, the least squares estimates of the individual regression coefficients ten to be **unstable** and can lead to erroneous inferences.
- In the last session we discussed the problem of multicollinearity and ways to diagnose this problem. We found that eliminating predictors from the analysis does not always work and in most analytical settings is not a feasible option.
- We consider two alternative approaches for dealing with multicollinearity:
  - Imposing or searching for constraints on the regression parameters.
  - Using alternative estimation techniques (e.g. principal components regression and ridge regression).

## Principal Components

- The principal components method is based on the fact that any set of $p$ predictors $X_1, X_2, \ldots, X_p$ can be **transformed** to a set of $p$ **orthogonal** variables.
- The new orthogonal variables are known as the **principal components** and are denoted by $C_1, C_2, \ldots, C_p$.
- Each variable $C_j$ is a linear function of the standardized variables $\tilde{X}_1, \tilde{X}_2, \ldots, \tilde{X}_p$.

$$C_j = v_{1j}\tilde{X}_1 + v_{2j}\tilde{X}_2 + \ldots + v_{pj}\tilde{X}_p \quad \text{for} \quad j = 1, 2, \ldots, p$$

- The coefficients of the linear functions are chosen so that the variables $C_1, \ldots, C_p$ are orthogonal.
- The coefficients for the $j$-th principal components $C_j$ are the elements of the $j$-th eigenvector that corresponds to the eigenvalue $\lambda_j$, the $j$-th largest eigenvalue of the correlation matrix of the $p$ variables.

$$V = \begin{pmatrix} V_1 & V_2 & \cdots & V_p \end{pmatrix} = \begin{pmatrix} v_{11} & v_{12} & \cdots & v_{1p} \\ v_{21} & v_{22} & \cdots & v_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ v_{p1} & v_{p2} & \cdots & v_{pp} \end{pmatrix}$$

# Example: French Econmony Data

P241

```
##    YEAR IMPORT DOPROD STOCK CONSUM
## 1   49   15.9  149.3   4.2  108.1
## 2   50   16.4  161.2   4.1  114.8
## 3   51   19.0  171.5   3.1  123.2
## 4   52   19.1  175.5   3.1  126.9
## 5   53   18.8  180.8   1.1  132.1
## 6   54   20.4  190.7   2.2  137.7
## 7   55   22.7  202.1   2.1  146.0
## 8   56   26.5  212.4   5.6  154.1
## 9   57   28.1  226.1   5.0  162.3
## 10  58   27.6  231.9   5.1  164.3
## 11  59   26.3  239.0   0.7  167.6
## 12  60   31.1  258.0   5.6  176.8
## 13  61   33.3  269.8   3.9  186.6
## 14  62   37.0  288.4   3.1  199.7
## 15  63   43.3  304.5   4.6  213.9
## 16  64   49.0  323.4   7.0  223.8
## 17  65   50.3  336.8   1.2  232.0
## 18  66   56.6  353.9   4.5  242.9
```

## Data Description

YEAR  Year of Observation.

IMPORT  Import Volume.

DOPROD  Domestic Production.

STOCK  Stock Formation.

CONSUM  Domestic Consumption.

Variables are measured in billion French francs.

## Principal Components

- It can be shown that the variance of the $j$-th principal component is $Var(C_j) = \lambda_j$ for $j = 1, 2, \ldots, p$. Therefore the variance-covariance matrix of the principal components is

$$\begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_p \end{pmatrix}$$

- All the off-diagonal elements are zero because the principal components are orthogonal. The value of the $j$-th diagonal element $\lambda_j$ is the variance of $C_j$, the $j$-th principal component.

- The principal components are arranged so that $\lambda_1 \geq \lambda_2 \geq \ldots \lambda_p$, which means that the first component has the largest variance.

## Principal Components

```
d <- head(P241[ ,c("DOPROD", "STOCK", "CONSUM")], 11)
d.pca <- prcomp(d, center=TRUE, scale=TRUE)
C <- d.pca$x
round(C, 4)
```

```
##         PC1     PC2     PC3
## 1   -2.1259  0.6387 -0.0207
## 2   -1.6189  0.5555 -0.0711
## 3   -1.1152 -0.0730 -0.0217
## 4   -0.8943 -0.0824  0.0108
## 5   -0.6442 -1.3067  0.0726
## 6   -0.1904 -0.6591  0.0266
## 7    0.3596 -0.7437  0.0428
## 8    0.9718  1.3541  0.0629
## 9    1.5593  0.9640  0.0236
## 10   1.7670  1.0152 -0.0450
## 11   1.9311 -1.6627 -0.0806
```

# Principal Components

```r
cormat <- cor(d)
eigen(cormat)    # Eigen Decomposition of Correlation Matrix


## eigen() decomposition
## $values
## [1] 1.999154934 0.998154176 0.002690889
##
## $vectors
##              [,1]         [,2]          [,3]
## [1,] 0.70633041  0.03568867   0.706982083
## [2,] 0.04350059 -0.99902908   0.006970795
## [3,] 0.70654444  0.02583046  -0.707197102
```

```r
round(var(C),4) # Variance-Covariance Matrix of PCs


##         PC1    PC2    PC3
## PC1 1.9992 0.0000 0.0000
## PC2 0.0000 0.9982 0.0000
## PC3 0.0000 0.0000 0.0027
```

## Remember

Multicollinearity leads to heterogeneous sizes of eigenvalues so that one eigenvalue is much smaller than the others. When one eigenvalue is exactly zero a perfect linear relationship (special case of extreme multicollinearity) among the original variables exists.

The variance-covariance matrix of the new variables only has entries on the main diagonal (which correspond to the eigenvalues) and zeros in all other places (as the variables are orthogonal).

## Principal Components

- The principal components lack simple interpretation as they are *a mixture* of the (standardized) original variables.
- Since $\lambda_j$ is the variance of the *j*-th principal component, a value of *lambda_j* $\approx 0$ shows that the respective principal componennt $C_j$ is equal to a constant. That constant is the mean value of $C_j$ (which is zero as the variables have been standardized).
- Inspecting the eigenvectors of the previous example shows that *only* the variables `CONSUM` and `DOPROD` play a relevant role when determining $C_3$.

```
##     X1_tilde      X2_tilde      X3_tilde
##   0.706982083  0.006970795  -0.707197102
```

$$\tilde{X}_1 \approx \tilde{X}_3 \quad \text{as} \quad v_{23} \approx 0.007 \approx 0$$

# Outline

## Principal Component Regression

- We consider the model for the *French Economomy Dataset*

$$\text{IMPORT} = \beta_0 + \beta_1\text{DOPROD} + \beta_2\text{STOCK} + \beta_3\text{CONSUM} + \epsilon$$

- This model expressed using the standardized variables $\tilde{Y} = (y_i - \bar{y})/s_y$ and $\tilde{X}_j = (x_{ij} - \bar{x}_j)/s_{x_j}$ yields

$$\tilde{Y} = \theta_1\tilde{X}_1 + \theta_2\tilde{X}_2 + \theta_3\tilde{X}_3 + \epsilon'$$

- Utilizing the principal components of the standardized predictors the model can be written as

$$\tilde{Y} = \alpha_1 C_1 + \alpha C_2 + \alpha_3 C_3 + \epsilon'$$

# Principal Component Regression

```r
# Data Preparation
d <- head(P241,11)
d_scaled <- as.data.frame(scale(d))
d_prcomp <- as.data.frame(cbind(IMPORT=d_scaled$IMPORT,
                          prcomp(d[,c("DOPROD","STOCK","CONSUM")],
                                 center=TRUE, scale=T)$x))
# Motel Estimation
mod1 <- lm(IMPORT ~  1 + DOPROD + STOCK + CONSUM, data=d)
(mod2 <- lm(IMPORT ~ -1 + DOPROD + STOCK + CONSUM, data=d_scaled))
```

```
##
## Call:
## lm(formula = IMPORT ~ -1 + DOPROD + STOCK + CONSUM, data = d_scaled)
##
## Coefficients:
##  DOPROD    STOCK   CONSUM
## -0.3393   0.2130   1.3027
```

```r
(mod3 <- lm(IMPORT ~ -1 + PC1 + PC2 + PC3, data=d_prcomp))
```

```
##
## Call:
## lm(formula = IMPORT ~ -1 + PC1 + PC2 + PC3, data = d_prcomp)
##
## Coefficients:
##    PC1     PC2     PC3
## 0.6900  0.1913  1.1597
```

## Principal Component Regression

```r
ev <- eigen(cor(d[,c("DOPROD","STOCK","CONSUM")]))$vectors

# Multiply eigenvectors with constant to match output in book
ev[,2]  <- ev[,2] * -1; ev[,3]  <- ev[,3] * -1

# Eigenvectors
ev
```

```
##            [,1]        [,2]         [,3]
## [1,] 0.70633041 -0.03568867 -0.706982083
## [2,] 0.04350059  0.99902908 -0.006970795
## [3,] 0.70654444 -0.02583046  0.707197102
```

## Principal Component Regression

- The coefficients of the principal component regression can be calculated based on the regression coefficients from the model using the standardized values.

$$
\begin{aligned}
\alpha_1 &= \phantom{-}0.706\ \theta_1 + \phantom{-}0.044\ \theta_2 + \phantom{-}0.707\ \theta_3 \\
\alpha_2 &= -0.036\ \theta_1 + \phantom{-}0.999\ \theta_2 + -0.026\ \theta_3 \\
\alpha_3 &= -0.707\ \theta_1 + -0.007\ \theta_2 + \phantom{-}0.707\ \theta_3
\end{aligned}
$$

- Conversely this relationship can be turned around to obtain the coefficients from the regression with standardized variables from the principal component regression.

$$
\begin{aligned}
\theta_1 &= 0.706\ \alpha_1 + -0.036\ \alpha_2 + -0.707\ \alpha_3 \\
\theta_2 &= 0.044\ \alpha_1 + \phantom{-}0.999\ \alpha_2 + -0.007\ \alpha_3 \\
\theta_3 &= 0.707\ \alpha_1 + -0.026\ \alpha_2 + \phantom{-}0.707\ \alpha_3
\end{aligned}
$$

# Principal Component Regression

```r
# Calculate alpha (principal components) from theta (standardized variables)
as.vector(coef(mod2) %*% ev)
```

```
## [1] 0.6899821 0.1913034 1.1596766
```

```r
coef(mod3)
```

```
##       PC1       PC2       PC3
## 0.6899821 0.1913034 1.1596766
```

```r
# Calculate theta (standardized variables) from alpha (principal components)
as.vector(ev %*% coef(mod3))
```

```
## [1] -0.3393426  0.2130484  1.3026815
```

```r
coef(mod2)
```

```
##      DOPROD      STOCK     CONSUM
## -0.3393426  0.2130484  1.3026815
```

# Principal Component Regression

$$\tilde{Y} = \theta_1 \tilde{X}_1 + \theta_2 \tilde{X}_2 + \theta_3 \tilde{X}_3 + \epsilon'$$
$$= \alpha_1 C_1 + \alpha C_2 + \alpha_3 C_3 + \epsilon'$$

- Although the above equations both hold, the $C$'s are **orthogonal**.
- The orthogonality bypasses (but not eliminates) the multicollinearity problem, however, the resulting relationship and therefore the coefficients are **not easily interpreted**.
- The $\alpha$'s unlike the $\theta$'s do not have simple interpretations as marginal effects of the original (standardized) predictor variables.

The final estimation results are always restated in terms of the $\theta$'s or origninal $\beta$'s for interpretation!

Based on the coefficients obtained from regressing the standardized variables the relationship can be expressed in terms of the original $\beta_j$'s using the following relationship:

$$\hat{\beta}_j = \frac{s_y}{s_j}\hat{\theta}_j \quad \text{for} \quad j = 1, 2, \ldots, p$$

$$\beta_0 = \bar{y} - \hat{\beta}_1 \bar{x}_1 - \hat{\beta}_2 \bar{x}_2 - \ldots - \hat{\beta}_p \bar{x}_p$$

This back-transform of the variables to the original scale is crucial for interpretation of the final results!

- Principal component regression can be used to **reduce collinearity in the estimation data**.
- this can be achieved by using **less than the full set of principal components** to explain the variation in the response.
- When all principal components are used the OLS solution can be exactly reproduced (as seen before).

- The $C_j$'s have sample variances $\lambda_1, \lambda_2, \ldots, \lambda_p$ equal to their eigenvalues.

```r
eigen(cor(d[,c("DOPROD","STOCK","CONSUM")]))$values
```

```
## [1] 1.999154934 0.998154176 0.002690889
```

- Since $C_3$ has very small variance, the linear function defining $C_3$ is **approximately equal to zero** and is the source of collinearity in the data.

## Reduction of Multicollinearity in the Data

■ We exclude $C_3$ from the analysis and consider the two possible remaining regression models

$$\tilde{Y} = \alpha_1 C_1 + \epsilon$$
$$\tilde{Y} = \alpha_1 C_1 + \alpha_2 C_2 + \epsilon$$

**There are two important things to note here:**

**1** In an regression equation where the full set of potential predictor variables under consideration are orthogonal, the estimated values of the regression **coefficients are not altered** when subsets of these variables are either introduced or deleted.

**2** Both models lead to estimates for **all** three of the original standardized coefficients $\theta_1, \theta_2$ and $theta_3$.

# Reduction of Multicollinearity in the Data

```
mod_prcomp1 <- lm(IMPORT ~ -1 + PC1            , data=d_prcomp)
mod_prcomp2 <- lm(IMPORT ~ -1 + PC1 + PC2      , data=d_prcomp)
mod_prcomp3 <- lm(IMPORT ~ -1 + PC1 + PC2 + PC3, data=d_prcomp)
```

|                      | Model 1    | Model 2    | Model 3    |
| -------------------- | ---------- | ---------- | ---------- |
| PC1                  | 0.69***    | 0.69***    | 0.69***    |
|                      | (0.05)     | (0.03)     | (0.02)     |
| PC2                  |            | 0.19***    | 0.19***    |
|                      |            | (0.04)     | (0.03)     |
| PC3                  |            |            | 1.16       |
|                      |            |            | (0.61)     |
| $R^2$                | 0.95       | 0.99       | 0.99       |
| Adj. $R^2$           | 0.95       | 0.99       | 0.99       |
| Num. obs.            | 11         | 11         | 11         |

$^{***}p < 0.001$; $^{**}p < 0.01$; $^{*}p < 0.05$

**Table 2:** Statistical models

# Reduction of Multicollinearity in the Data

```r
# Coefficients for standardized predictors when using one principal component
coefs1 <- coef(mod3)[1] * ev[,1]
names(coefs1) <- c("DOPROD","STOCK","CONSUM")
coefs1
```

```
##    DOPROD      STOCK     CONSUM
## 0.48735534 0.03001463 0.48750301
```

```r
# Coefficients for standardized predictors when using two principal components
coefs2 <- coef(mod3)[1] * ev[,1] + coef(mod3)[2] * ev[,2]
names(coefs2) <- c("DOPROD","STOCK","CONSUM")
coefs2
```

```
##    DOPROD     STOCK    CONSUM
## 0.4805280 0.2211323 0.4825616
```

# Reduction of Multicollinearity in the Data

```
s <- apply(d[ ,c("IMPORT","DOPROD","STOCK","CONSUM")],2,sd)
m <- apply(d[ ,c("IMPORT","DOPROD","STOCK","CONSUM")],2,mean)
```

```
# Model with one PC for non-standardized data
coefs_org1 <- s[1]/s[2:4] * coefs1
intercept_org1 <- unname(m[1] - sum(m[2:4]*coefs_org1))
(beta_org1 <- c(Intercept=intercept_org1, coefs_org1))
```

```
##    Intercept       DOPROD        STOCK       CONSUM
## -7.74582557   0.07381387   0.08269039   0.10734749
```

```
# Model with two PCs for non-standardized data
coefs_org2 <- s[1]/s[2:4] * coefs2
intercept_org2 <- unname(m[1] - sum(m[2:4]*coefs_org2))
(beta_org2 <- c(Intercept=intercept_org2, coefs_org2))
```

```
##    Intercept       DOPROD        STOCK       CONSUM
## -9.13010782   0.07277981   0.60922012   0.10625939
```

## Reduction of Multicollinearity in the Data

- The following table shows that the coefficients are dependent on the number of incorporated principal components.
- As each component explains additional variance the $R^2$ inceases with the number of considered principal components.

|           | std_PC1 | org_PC1 | std_PC2 | org_PC2 | std_PC3 | org_PC3 |
|-----------|---------|---------|---------|---------|---------|---------|
| Intercept | NA      | -7.746  | NA      | -9.130  | NA      | -10.128 |
| DOPROD    | 0.487   | 0.074   | 0.481   | 0.073   | -0.339  | -0.051  |
| STOCK     | 0.030   | 0.083   | 0.221   | 0.609   | 0.213   | 0.587   |
| CONSUM    | 0.488   | 0.107   | 0.483   | 0.106   | 1.303   | 0.287   |

# Outline