

# Multiple Linear Regression

# Examples

- 晶圓(Wafer)良率或不良率～多個機台+機台參數
- 螢幕Mura～機台+材料面板瑕疵
- 汽車保養頻率～引擎+組裝+動力系統
- 眼科疾病～年紀+糖尿病+血壓+家族病史
- GDP～人口數+出口+進口+...

# Multiple linear regression model

$$Y \sim \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \cdots$$

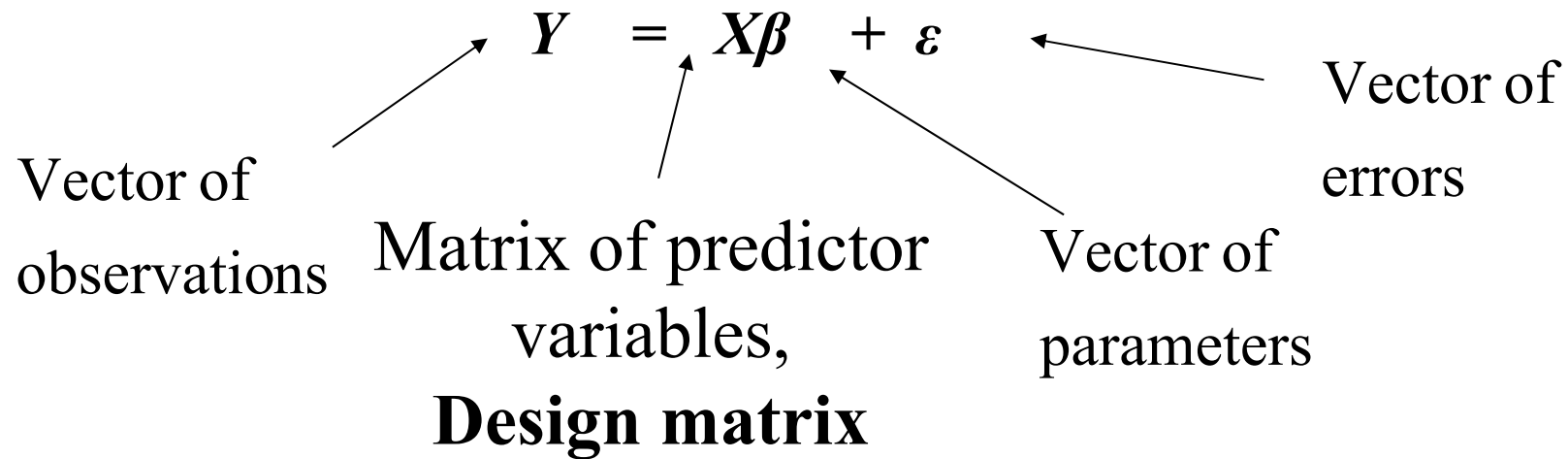
$$Y \sim \beta_0 + \beta_1 X_1 + \beta_2 X_2 = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2$$

(e.g.  $X_2 = X_1^2$ )

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i, \quad i = 1, \cdots, n$$

- Condition(A1).  $E(\varepsilon_i) = 0, \quad i = 1, \cdots, n.$
- Condition(A2).  $Var(\varepsilon_i) = \sigma^2, \quad i = 1, \cdots, n.$
- Condition(A3).  $Cov(\varepsilon_i, \varepsilon_j) = 0, \text{ for } i \neq j.$

# Regression in Matrix Terms



$$\mathbf{Y} = \underset{\sim}{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

## RSS (SSE) and LSE

$$\begin{aligned} RSS(\boldsymbol{\beta}) &= \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n \left( y_i - (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}) \right)^2 \\ &= \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^T \mathbf{y} - 2\boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{y}) + \boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X}) \boldsymbol{\beta} \end{aligned}$$

Derivatives the RSS and let it be zero, we have

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $E\left(\hat{\boldsymbol{\beta}} \mid \boldsymbol{X}\right) = \boldsymbol{\beta}$

- $Cov\left(\hat{\boldsymbol{\beta}} \mid \boldsymbol{X}\right) = \left(\boldsymbol{X}^T \boldsymbol{X}\right)^{-1} \sigma^2$

- $\hat{\boldsymbol{y}} = \boldsymbol{X} \hat{\boldsymbol{\beta}} = \boldsymbol{X} \left(\boldsymbol{X}^T \boldsymbol{X}\right)^{-1} \boldsymbol{X}^T \boldsymbol{y} = \boldsymbol{H} \boldsymbol{y}$

***H***: projection matrix or hat matrix

$$h_{ii} = \mathbf{x}_i^T \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_i^T$$

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I} - \mathbf{H}) \mathbf{y} = \mathbf{M}\mathbf{y}$$

Notes: ***H*** and ***M*** are symmetric and idempotent.



$$\hat{\boldsymbol{\varepsilon}}^T \boldsymbol{X} = \boldsymbol{0}$$

$$\hat{\boldsymbol{\varepsilon}}^T \boldsymbol{1} = 0$$

$$\hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{y}} = 0$$

$$\hat{\sigma}^2 = s^2 = \frac{\hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}}{n - (p + 1)}$$

$$SST(n-1) = SSR(p) + SSE(n - (p+1))$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\Leftrightarrow (\mathbf{y} - \mathbf{I}\bar{y})^T (\mathbf{y} - \mathbf{I}\bar{y}) = (\hat{\mathbf{y}} - \mathbf{I}\bar{y})^T (\hat{\mathbf{y}} - \mathbf{I}\bar{y}) + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}$$

$$\Leftrightarrow \mathbf{y}^T \left( \mathbf{I} - \frac{\mathbf{1}\mathbf{1}^T}{n} \right) \mathbf{y} = \mathbf{y}^T \mathbf{H}^T \left( \mathbf{I} - \frac{\mathbf{1}\mathbf{1}^T}{n} \right) \mathbf{H}\mathbf{y} + \mathbf{y}^T \mathbf{M}\mathbf{y}$$

$$\Leftrightarrow \mathbf{y}^T \mathbf{L}\mathbf{y} = \mathbf{y}^T \mathbf{H}^T \mathbf{L}\mathbf{H}\mathbf{y} + \mathbf{y}^T \mathbf{M}\mathbf{y}$$

$$\begin{aligned}
 R^2 &= \frac{SSR}{SST} = \frac{\text{Variability explained by the model}}{\text{Total sample variability}} = \frac{\mathbf{y}^T \mathbf{H}^T \mathbf{L} \mathbf{H} \mathbf{y}}{\mathbf{y}^T \mathbf{L} \mathbf{y}} \\
 &= 1 - \frac{\text{Unexplained(error)}}{\text{Total sample variability}} = 1 - \frac{\mathbf{y}^T \mathbf{M} \mathbf{y}}{\mathbf{y}^T \mathbf{L} \mathbf{y}} \\
 &= 1 - \frac{SSE}{SST}
 \end{aligned}$$

$$0 \leq R^2 \leq 1$$

# ANOVA table

$H_0 : \beta_1 = \cdots = \beta_p = 0$  versus  $H_1 : \beta_j \neq 0$  for some  $j$

Source of variation (Source)	Sum of squares (SS)	Degrees of freedom (df)	Mean square (MS)	F value (F)	P-value
Regression	SSR	p	$MSR = SSR / p$	$F = MSR / MSE$	$P(F_{p, n-(p+1)} > F)$
Error	SSE	n-(p+1)	$MSE = SSE / (n - (p + 1))$		
Total	SST	n-1			

## Testing whether a specified subset of the predictors have coefficients equal to 0

Recall:

$H_0 : \beta_1 = \dots = \beta_p = 0$  versus  $H_1 : \beta_j \neq 0$ , for some  $j = 1, \dots, p$

i.e.  $Y_i = \beta_0 + \varepsilon_i$ , for  $i = 1, \dots, n$  and

$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$ , for  $i = 1, \dots, n$

and for at least one of  $\beta_j \neq 0$ ,  $j = 1, \dots, p$

$$F = \frac{SSR / p}{SSE / (n - (p + 1))} \sim F(p, n - (p + 1))$$

Reject region  $C = \{F > F_\alpha(p, n - (p + 1))\}$

$H_0 : \beta_{m+1} = \dots = \beta_p = 0$  versus  $H_1 : \beta_j \neq 0$ , for some  $j = m+1, \dots, p$   
two nested models,  $H_0$  : reduced model versus  $H_1$  : full model

$H_0 : Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im} + \varepsilon_i$ , for  $i = 1, \dots, n$  and  
 $H_1 : Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im} + \beta_{m+1} x_{im+1} + \dots + \beta_p x_{ip} + \varepsilon_i$ ,  
for  $i = 1, \dots, n$  and for at least one of  $\beta_j \neq 0$ ,  $j = m+1, \dots, p$

$$(\text{partial}) F = \frac{\left( SSR_{n-(m+1)} - SSR_{n-(p+1)} \right) / (p-m)}{SSE / (n-(p+1))} \sim F(p-m, n-(p+1))$$

$$\text{Reject region } C = \left\{ F > F_\alpha(p-m, n-(p+1)) \right\}$$

$H_0 : \beta_j = 0$  versus  $H_1 : \beta_j \neq 0$ , for  $j \in \{1, \dots, p\}$

two nested models,  $H_0$  : reduced (restricted model) model versus

$H_1$  : full (unrestricted model) model

$$H_0 : Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{j-1} x_{i,j-1} + \beta_{j+1} x_{i,j+1} + \dots + \beta_p x_{ip} + \varepsilon_i,$$

for  $i = 1, \dots, n$ ,  $j \in \{1, \dots, p\}$  and

$$H_1 : Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i, \text{ for } i = 1, \dots, n$$

$$(\text{partial}) F = \frac{\left( SSR_{n-p} - SSR_{n-(p+1)} \right)}{SSE / (n - (p + 1))} \sim F(1, n - (p + 1))$$

$$\text{Reject region } C = \left\{ F > F_\alpha(1, n - (p + 1)) \right\}$$

# Sequential F test statistics

For  $Y_i = \beta_0 + \beta_1 x_{i1} + \varepsilon_i$ ,  $i = 1, \dots, n$ , consider

$H_0 : \beta_1 = 0$  versus  $H_1 : \beta_1 \neq 0$ .  $\Rightarrow F_{1,n-2}$

Next, for  $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i$ ,  $i = 1, \dots, n$ , consider

$H_0 : \beta_2 = 0$  versus  $H_1 : \beta_2 \neq 0$ .  $\Rightarrow F_{1,n-3}$

$\vdots$

For  $Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$ ,  $i = 1, \dots, n$ , consider

$H_0 : \beta_p = 0$  versus  $H_1 : \beta_p \neq 0$ .  $\Rightarrow F_{1,n-p-1}$

Then, the nested F test statistic,  $F_{1,n-2}, F_{1,n-3}, \dots, F_{1,n-p-1}$ , are called sequential F test statistic.



# Dummy variables

## **Binary variable**

- sex (female or male)
- disease (yes or no)
- master degree (yes or no)
- ...

# Parallel regression lines

$$Y_i = \beta_0 + \beta_1 x_i + \delta d_i + \varepsilon_i, \quad i = 1, \dots, n$$

$x_i$  : continuous covariate

$d_i$  : binary covariate

$$Y_i = \begin{cases} \beta_0 + \beta_1 x_i + \varepsilon_i & \text{if } d_i = 0 \\ (\beta_0 + \delta) + \beta_1 x_i + \varepsilon_i & \text{if } d_i = 1 \end{cases},$$

$d_i = 0$  : base group

# Regression lines with equal intercepts but different slopes

$$Y_i = \beta_0 + \beta_1 x_i + \delta' d_i x_i + \varepsilon_i, \quad i = 1, \dots, n$$

$x_i$  : continuous covariate

$d_i$  : binary covariate

$$Y_i = \begin{cases} \beta_0 + \beta_1 x_i + \varepsilon_i & \text{if } d_i = 0 \\ \beta_0 + (\beta_1 + \delta') x_i + \varepsilon_i & \text{if } d_i = 1 \end{cases},$$

$d_i = 0$  : base group

# Unrelated regression lines

$$Y_i = \beta_0 + \beta_1 x_i + \delta d_i + \delta' d_i x_i + \varepsilon_i, \quad i = 1, \dots, n$$

$x_i$  : continuous covariate

$d_i$  : binary covariate

$$Y_i = \begin{cases} \beta_0 + \beta_1 x_i + \varepsilon_i & \text{if } d_i = 0 \\ (\beta_0 + \delta) + (\beta_1 + \delta') x_i + \varepsilon_i & \text{if } d_i = 1 \end{cases},$$

$d_i = 0$  : base group

# Dummy variables (cont')

## **Categorical variable**

- salary level (150+, 100~150, 50~100, < 50)
- education level(phd, master, bachelor, other)
- ethnic ("Aethiopian or Black", "Caucasian or White", "Mongolian or Yellow", "American or Red" and "Malayan or Brown")
- ...

## Categorical variable: four levels

	d1	d2	d3
C1	0	0	0
C2	1	0	0
C3	0	1	0
C4	0	0	1

# Parallel regression lines

$$Y_i = \beta_0 + \beta_1 x_i + \delta_1 d_{i1} + \delta_2 d_{i2} + \delta_3 d_{i3} + \varepsilon_i, \quad i = 1, \dots, n$$

$x_i$  : continuous covariate

$d_i$  : dummy variables for the categorical covariate

$$Y_i = \begin{cases} \beta_0 + \beta_1 x_i + \varepsilon_i & \text{if C1} \\ (\beta_0 + \delta_1) + \beta_1 x_i + \varepsilon_i & \text{if C2} \\ (\beta_0 + \delta_2) + \beta_1 x_i + \varepsilon_i & \text{if C3} \\ (\beta_0 + \delta_3) + \beta_1 x_i + \varepsilon_i & \text{if C4} \end{cases},$$

$(d_1, d_2, d_3) = (0, 0, 0)$  : base group

# Categorical variable: four levels

	d1	d2	d3
C1	0	0	0
C2	1	0	0
C3	1	1	0
C4	1	1	1



# Collinearity in linear model with the intercept term

	d1	d2	d3	D4
C1	1	0	0	0
C2	0	1	0	0
C3	0	0	1	0
C4	0	0	0	1

# Multi-collinearity

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon.$$

What is multicollinearity?

If there is a set of  $c_1, \dots, c_p$  (not all zero) for which

$$\sum_{j=1}^p c_j x_j \cong 0.$$

## What is the impact of multicollinearity?

The LSE,  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ , will not be stable.

Consider the following special case,

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\text{Var}(\varepsilon_i) = \sigma^2$  and all the  $\varepsilon_i$  are uncorrelated.

Assume that  $\mathbf{X}_1$  and  $\mathbf{X}_2$  have been scaled so that

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2 \cdot \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}{1 - r_1^2} = \frac{\sigma^2}{1 - r_1^2}.$$

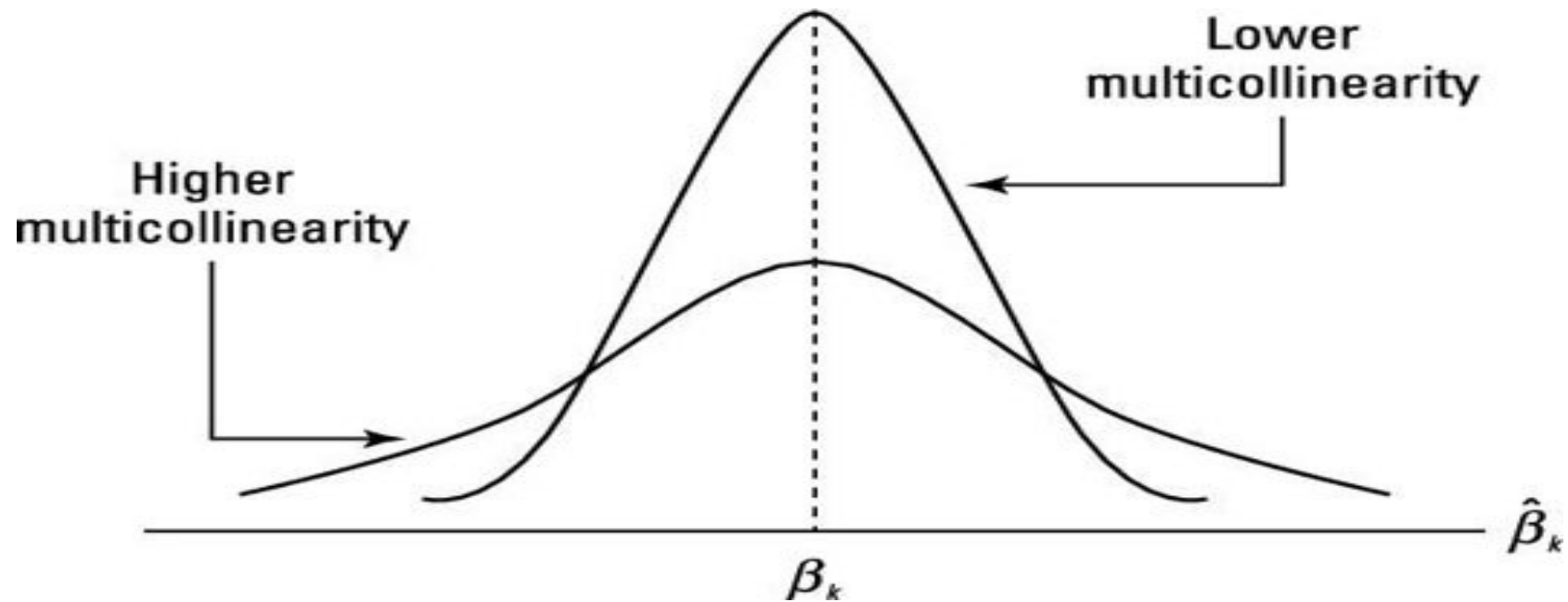
$1 - r_1^2$  is the residual mean square error (RMSE) obtained from regressing  $X_1$  and  $X_2$ . In particular,

$$r_1^2 = \frac{\left( \sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) \right)^2}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}.$$

- When  $r_1^2 = 0$ ,  $Var(\hat{\beta}_1)$  has its minimum value.
- As  $r_1^2 \rightarrow 1$ ,  $Var(\hat{\beta}_1)$  will get to be arbitrarily large.

# Variance inflation factors (VIFs)

$$\text{VIFs} = \frac{1}{1 - r_i^2}, i = 1, \dots, p$$
 represent the inflation that each regression expression above ideal.



# The condition number

It can be measured in terms of the ratio of the largest to the smallest eigenvalue, e.g., the quantity

$$\phi = \frac{\lambda_{\max}}{\lambda_{\min}}$$

which is called the condition number of the correlation matrix.

Ratios of eigenvalues, i.e., ratios  $\phi_j = \lambda_{\max} / \lambda_j$  are more reliable for diagnosing the impact of a dependency than the eigenvalue  $\lambda_j$  itself.

Consider centered and scaled regressors,  $\mathbf{X}^*$ ,  $\mathbf{X}^{*T} \mathbf{X}^*$  being a correlation matrix and

$$\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^{*T} \mathbf{X}^*)^{-1}, \quad \sum_{i=1}^p \frac{\text{Var}(\hat{\beta}_i)}{\sigma^2} = \text{tr}(\mathbf{X}^{*T} \mathbf{X}^*)^{-1} = \sum_{i=1}^p \frac{1}{\lambda_i}.$$

The operation given by  $\mathbf{V}^T (\mathbf{X}^{*T} \mathbf{X}^*)^{-1} \mathbf{V} = \text{diag}(\lambda_1, \dots, \lambda_p)$

$\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_p)$  is an orthogonal matrix, is called eigenvalue decomposition of  $\mathbf{X}^{*T} \mathbf{X}^*$ .

- If multicollinearity is present, at least one value of  $j$ ,

$\mathbf{v}_j^T (\mathbf{X}^{*T} \mathbf{X}^*)^{-1} \mathbf{v}_j (\geq \lambda_{\min}) \cong 0$ , this implies that for at least one eigenvector  $\mathbf{v}_j$ ,

$$\sum_{\ell=1}^p \mathbf{v}_\ell \mathbf{x}_j^* \cong \mathbf{0}.$$

# VIF in R

- <https://www.rdocumentation.org/packages/car/versions/3.0-7/topics/vif>
- `(install.packages("car")); library(car);`

```
set.seed(1082); n = 100 ;
```

```
x1 = rgamma(n, shape = 1) ; x2 = rexp(n) ; y = 3+ 2*x1 + x2 + rnorm(n)
```

```
vif(lm(y~x1+x2)) ## vifs < 1.1
```

```
set.seed(1082); n = 100 ;
```

```
x1 = rgamma(n, shape = 1) ; x2 = 2*x1 + rnorm(n, sd = 0.5) ; y = 3+ 2*x1 + x2 +  
rnorm(n)
```

```
vif(lm(y~x1+x2)) ## vifs > 20
```



# Alternatives to least squares in cases of multicollinearity

- Ridge regression
- Principal components regression
- Variable selection

# Ridge regression

- The ridge regression estimator of the coefficient  $\boldsymbol{\beta}$  is found by solving for  $\boldsymbol{\beta}_R$  in the system of equations

$(\mathbf{X}^T \mathbf{X} + \lambda_R \mathbf{I}) \boldsymbol{\beta}_R = \mathbf{X}^T \mathbf{y}$  where  $\lambda_R \geq 0$  often referred to as a shrinkage parameter. The solution is given by

$$\hat{\boldsymbol{\beta}}_R = (\mathbf{X}^T \mathbf{X} + \lambda_R \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}.$$

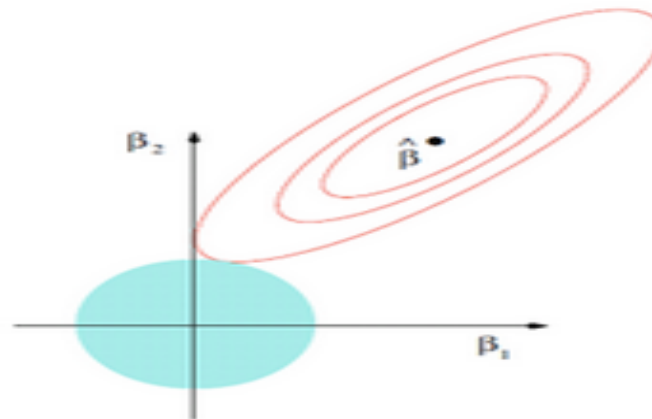
- Consider centered and scaled regressors,

$$\sum_{i=1}^p \frac{\text{Var}(\hat{\beta}_{i,R})}{\sigma^2} = \sum_{i=1}^p \frac{\lambda_i}{(\lambda_i + \lambda_R)^2}, \lambda_i, i = 1, \dots, p \text{ are the eigenvalue of } \mathbf{X}^T \mathbf{X}.$$

## Ridge Regression:

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2,$$

subject to  $\sum_{j=1}^p \beta_j^2 \leq t, \quad (\text{L2 term})$



[https://rpubs.com/skydome20/R-Note18-Subsets\\_Shrinkage\\_Methods](https://rpubs.com/skydome20/R-Note18-Subsets_Shrinkage_Methods)

- For example, in the case of the three regressor variables with  $\lambda = (2.985, 0.01, 0.005)$ , least squares estimation gives

$$\sum_{i=1}^3 \frac{\text{Var}(\hat{\beta}_i)}{\sigma^2} = \sum_{i=1}^p \frac{1}{\lambda_i} = 300.335$$

- If ridge regression with say  $\lambda_R = 0.1$  is used, the sum of the variances is given by

$$\sum_{i=1}^3 \frac{\lambda_i}{(\lambda_i + \lambda_R)^2} \cong 1.6$$

- The procedure would be successful if a  $\lambda_R$  is chosen so that the variance reduction is greater than the bias term given in

$$\sum_{i=1}^p \left( E(\hat{\beta}_{i,R}) - \beta_i \right)^2 = \lambda_R^2 \boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X} + \lambda_R \mathbf{I})^{-2} \boldsymbol{\beta}.$$

# Shrinkage

- the similarity to the ordinary least squares solution

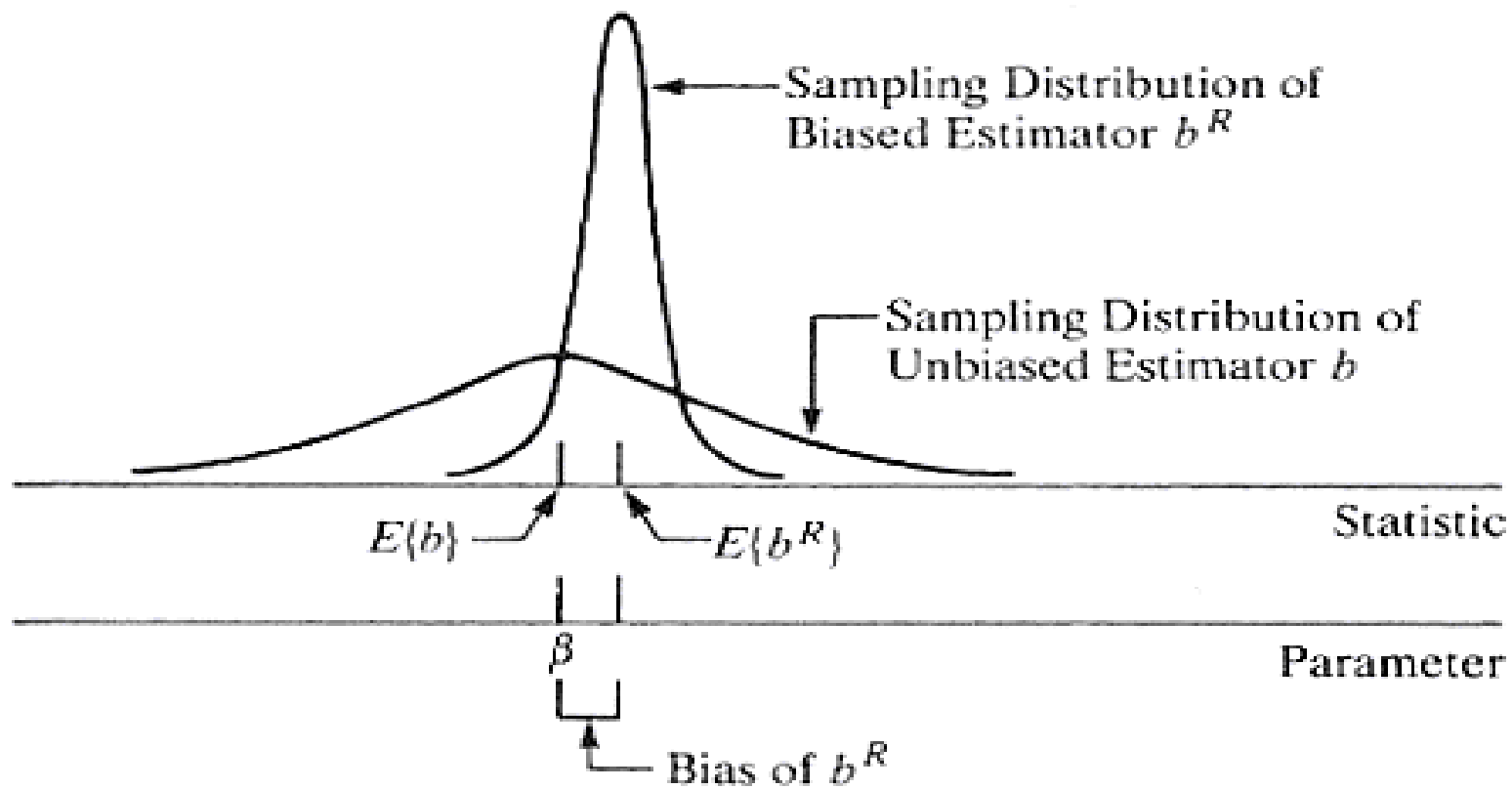
i.e. As  $\lambda_R \rightarrow \infty$ ,  $\hat{\beta}_R \rightarrow \hat{\beta}_{OLS}$ .

- In the special case of an orthonormal design matrix,

$$\hat{\beta}_R = \frac{\hat{\beta}_{OLS}}{1 + \lambda_R}$$

- This illustrates the essential feature of ridge regression : shrinkage.
- Applying the ridge regression penalty has the effect of shrinking the estimates toward zero - introducing bias but reducing the variance of the estimate

# Biased estimator with small variance



# Ridge regression in R

- <https://stat.ethz.ch/R-manual/R-devel/library/MASS/html/lm.ridge.html>
- `install.packages("MASS")`
- `install.packages("glmnet")`

# Principal components regression (PCA)

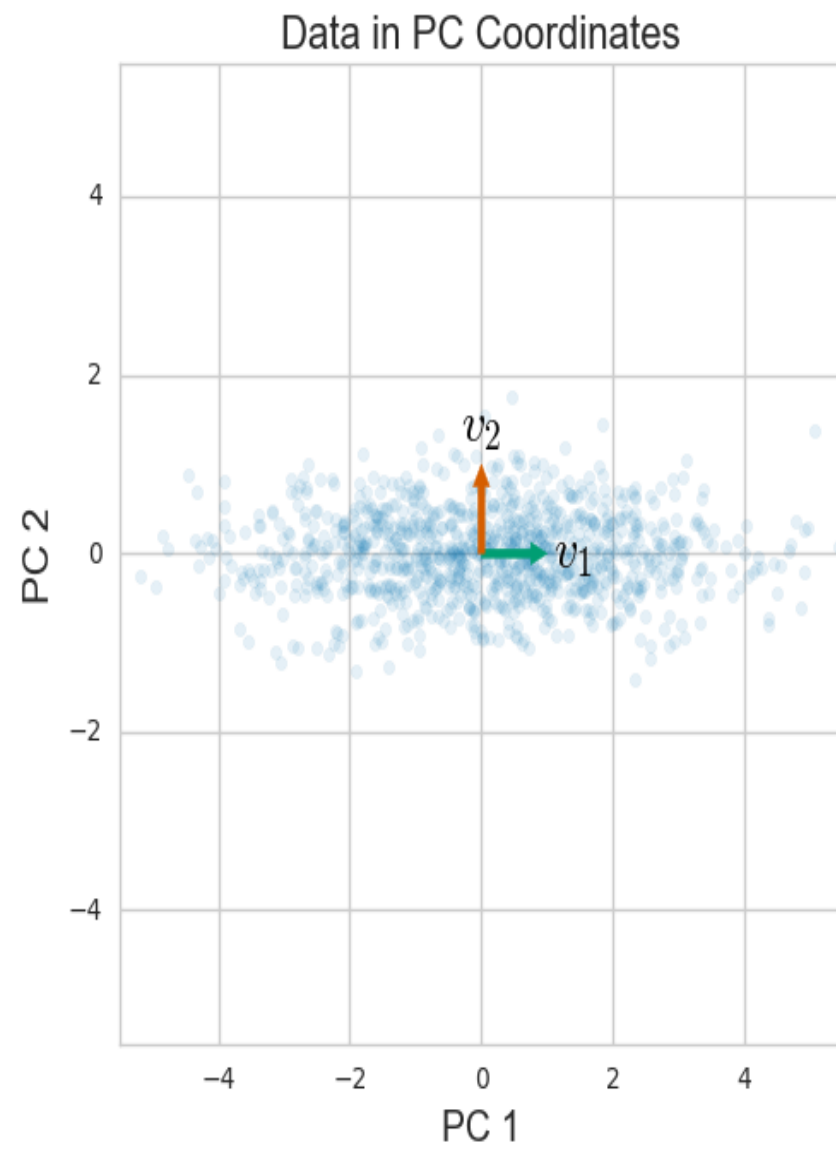
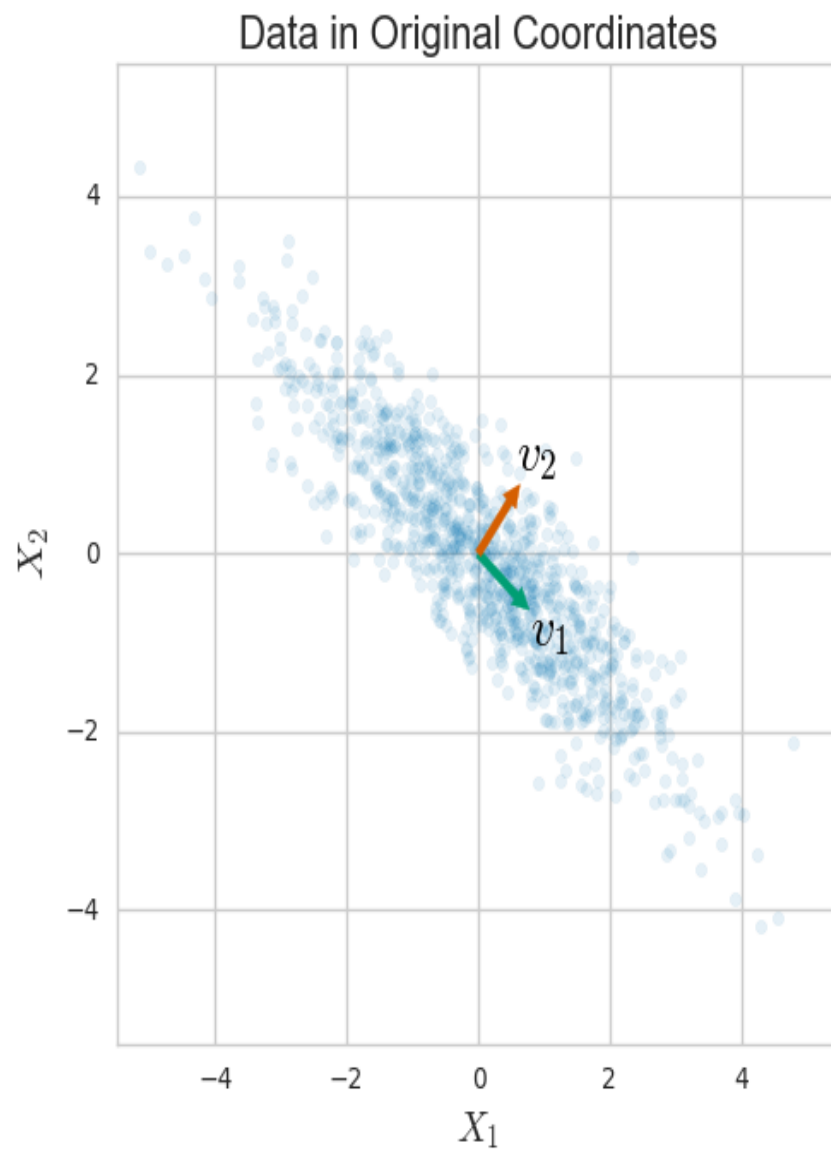
- Consider the original regression model in the form

$$\mathbf{y} = \beta_0 \mathbf{1} + \mathbf{X}^* \mathbf{V} \mathbf{V}^T \boldsymbol{\beta} + \boldsymbol{\varepsilon} = \beta_0 \mathbf{1} + \mathbf{Z} \boldsymbol{\alpha} + \boldsymbol{\varepsilon}. \text{ We have}$$

$\mathbf{Z}^T \mathbf{Z} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$  where  $\mathbf{V} \mathbf{V}^T = \mathbf{I}$  and  $\mathbf{V}$  is an orthogonal matrix.

- $\frac{\text{Var}(\hat{\alpha}_i)}{\sigma^2} = \frac{1}{\lambda_i}, i = 1, \dots, p.$





- Suppose we consider the matrix  $V = [v_1, \dots, v_p]$  of normalized eigenvectors of  $X^{*T} X^*$  partitioned into  $V = [V_k : V_r]$  and consider  $\Lambda$  to be a diagonal matrix of eigenvalues of  $X^{*T} X^*$ . We partition  $\Lambda$  as  $\Lambda = [\text{diag}(\Lambda_k) : \text{diag}(\Lambda_r)]$ .
- An elimination of one (at least one) principal component, that associated with the small eigenvalue, may substantially reduce the total variance in the model and thus produce an appreciably improved prediction equation.
- The difficulty arises in the decision of how many components we should eliminate.

## Advantages on principal component regression (PCR)

- There is virtually no limit for the number of predictors.  
i.e.  $p$  is allowed to be larger than  $n$
- Correlated or noisy predictors do not undermine the regression fit.
- The PCs carry the maximum amount of variance possible.

# Partial least squares (PLS) regression

- Stepwise **select** the orthogonal variables of PCA into the model and coefficients are estimated by **ordinary least squares**
- An alternative to PLS, principal component regression (PCR) i.e. fit using the first  $k$  principal components from the predictors  
drawback: no guarantee that PCs are associated with the outcome

# PCA and PCR in R

- <https://www.cnblogs.com/leezx/p/6120302.html>
- <https://www.jamleecute.com/principal-components-analysis-pca-%E4%B8%BB%E6%88%90%E4%BB%BD%E5%88%86%E6%9E%90/>
- <https://cran.r-project.org/web/packages/pls/pls.pdf>

```
> set.seed(1082); n = 100 ; x1 = rgamma(n, shape = 1) ; x2 = 2*x1 + rnorm(n, sd = 0.25) ; x3 = rexp( n )
> x.input = cbind(x1,x2,x3) ; prcomp( x.input )
Standard deviations (1, .., p=3):
[1] 2.33396740 0.92529449 0.09974213
```

```
Rotation (n x k) = (3 x 3):
      PC1      PC2      PC3
x1  0.44365914 -0.007453146  0.89616461
x2  0.89479755 -0.052149124 -0.44341607
x3 -0.05003904 -0.998611496  0.01646738
> biplot( prcomp(x.input), scale = 0) ; summary( prcomp(x.input))
Importance of components:
      PC1      PC2      PC3
Standard deviation  2.3340 0.9253 0.09974
Proportion of Variance 0.8628 0.1356 0.00158
Cumulative Proportion 0.8628 0.9984 1.00000
```

# $R^2$ and variable selection

A computing formula for the square of the multiple coefficient coefficient in a  $k$  – parameter model,  $J_k$ , is

$$R_{J_K}^2 = \frac{SSR_{J_K}}{SST} = 1 - \frac{SSE_{J_K}}{SST} = 1 - \frac{\hat{\sigma}_{J_K}^2}{\hat{\sigma}^2}.$$