

①

Partial Least Squares regression

Let $X = (X_1, X_2, \dots, X_p)$ be
a random vector with values in \mathbb{R}^p .

and

$Y = (Y_1, \dots, Y_q)$ a r.v. $Y \in \mathbb{R}^q$

X : predictor

Y : response.

Regression problem :

Approximate : $\mathbb{E}(Y|X)$
(conditional expectation).

②

$$E(Y|X) = f(X), \quad f: \mathbb{R}^p \rightarrow \mathbb{R}^2$$

$$f = (f_1, \dots, f_2), \quad f_i: \mathbb{R}^p \rightarrow \mathbb{R}.$$

Linear model:

$$f(x) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p,$$

with $\beta_i \in \mathbb{R}^2, i=0, \dots, p.$



$$\begin{cases} f_1(x) = \beta_0^1 + \beta_1^1 X_1 + \dots + \beta_p^1 X_p \\ \vdots \\ f_2(x) = \beta_0^2 + \beta_1^2 X_1 + \dots + \beta_p^2 X_p \end{cases}$$

Put $\beta = (\beta_0, \beta_1, \dots, \beta_p)$

③

Least Squares estimator:

$$\hat{\beta} = \arg \min_{\beta} \mathbb{E} \left(\left\| Y - \langle \tilde{X}, \beta \rangle_{\mathbb{R}^{q+1}} \right\|_{\mathbb{R}^q}^2 \right)$$

where $\tilde{X} = (1, X_1, \dots, X_p)$.

and

$$\langle \tilde{X}, \beta \rangle_{\mathbb{R}^{q+1}} = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

$$\hat{\beta} = \arg \min_{\beta} \mathbb{E} \left(\sum_{i=1}^q (Y_i - f_i(x))^2 \right)$$

But

$$\min \mathbb{E} \left(\sum_{i=1}^q (Y_i - f_i(x))^2 \right) = \sum_{i=1}^q \min \mathbb{E} (Y_i - f_i(x))^2$$

\Leftrightarrow q independent minimization problems
($q = 1$).

④

Remark : The estimation of the linear model by the least squares method does not take into account the (link) dependence between Y_i 's!

⑤ Let us consider $q = 1$ ($Y \in \mathbb{R}$)
and $\mathbb{E}(Y) = 0$
 $\mathbb{E}(X_i) = 0, i = 1, \dots, p$ } $\Leftrightarrow \beta_0 = 0.$

Then

$$\hat{\beta} = V^{-1} \cdot \mathbb{E}(X \cdot Y) \quad \text{where}$$

$$V = \{v_{ij}\}_{1 \leq i, j \leq p}, \quad v_{ij} = \mathbb{E}(X_i \cdot X_j)$$

and

$$\mathbb{E}(X \cdot Y) = \begin{pmatrix} \mathbb{E}(X_1 \cdot Y) \\ \vdots \\ \mathbb{E}(X_p \cdot Y) \end{pmatrix} \in \mathbb{R}^p$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix} \in \mathbb{R}^p$$

⑥

Denote by

$$\hat{Y} = X_1 \hat{\beta}_1 + \dots + X_p \hat{\beta}_p$$

$$= X^T V^{-1} E(X \cdot Y)$$

the least square estimation of the conditional expectation in the linear model.

Remark

$P_x = X^T V^{-1} E(X \cdot)$ is the orthogonal projector on the linear space spanned by $\{X_1, \dots, X_p\}$ with the dot product $\langle X_1, X_2 \rangle = E(X_1 X_2)$

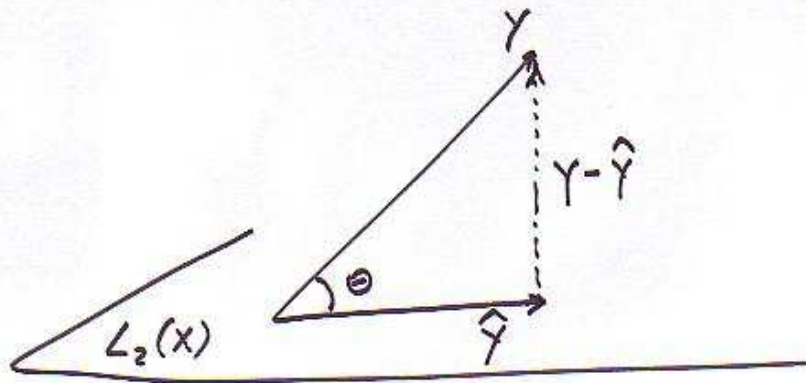
$$P_x(z) = X^T V^{-1} E(X \cdot z) \quad \forall z \in L_2(\Omega)$$

and

$$P_x^2(z) = P_x(z).$$

Thus $\hat{Y} = P_x(Y).$

(7)



$$\varepsilon = y - \hat{y} \quad (\text{error})$$

Analysis of variance :

$$V(y) = V(\hat{y}) + V(y - \hat{y})$$

(1)

and since $E(y) = 0$ and $E(x_i) = 0$:

$$E(y^2) = E(\hat{y}^2) + \underline{E(y - \hat{y})^2}$$

↑
objective function
of the least squares
criterion!

(2)

$$E(y - \hat{y})^2 = (1 - R^2) E(y^2) \quad (R^2 = \cos^2(\theta)).$$

⑧

$$E(Y - \hat{Y})^2 = (1 - \underline{R^2}) E(Y^2)$$

where R^2 is the correlation coefficient
between Y and $\hat{Y} = \beta_0 X_0 + \dots + \beta_p X_p$.

Thus minimize the least squares criterion
is equivalent to maximize the value
of R^2 !

⑨

Estimation from a sample of size n

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad X = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \dots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{np} \end{pmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad \text{and}$$

$$\hat{Y} = X (X^T X)^{-1} X^T Y = X \hat{\beta}.$$

Remark if data is not centered
then:

$$\hat{\beta}_0 = \bar{Y} - (\hat{\beta}_1 \bar{X}_1 + \dots + \hat{\beta}_p \bar{X}_p) \quad \text{with}$$

$$\bar{Y} = \frac{Y_1 + \dots + Y_n}{n}$$

$$\bar{X}_i = \frac{X_{1i} + \dots + X_{ni}}{n}, \quad i=1, \dots, p.$$

⑩

Problems of the linear model estimation

- Multicollinearity of the predictor

$\exists \alpha_1, \alpha_2, \dots, \alpha_p$ constants

$$\underline{\alpha_1 X_1 + \dots + \alpha_p X_p = 0} \quad \text{p.s.}$$

It occurs

- by the nature of the X_i 's

- when $n \leq p$

Consequence : V^{-1} does not exist !

\Rightarrow Bad estimation of β_i 's.

11

$$\boxed{V(\hat{\beta}) = \frac{\sigma^2}{n} \cdot V^{-1}}$$

where $\sigma^2 = E(Y - \hat{Y})^2 = \text{Var}(Y - \hat{Y})$
↑ residual variance.

$V(\hat{\beta})$ is then estimated by

$$\boxed{V(\hat{\beta}) \approx \hat{\sigma}^2 \cdot (X^T X)^{-1}}$$
 with

$$\hat{\sigma}^2 = \frac{1}{n-p-1} \cdot \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

If multicollinearity:

$V(\hat{\beta}_i) \rightarrow \infty$ when $\det(X^T X) \rightarrow 0$

\Rightarrow student test for $\begin{cases} H_0: \beta_i = 0 \\ H_a: \beta_i \neq 0 \end{cases}$

is not significant!

(12)

Remark The student test for

$$\begin{cases} H_0 : \beta_i = 0 \\ H_1 : \beta_i \neq 0 \end{cases}$$

is a test for the contribution of the variable X_i , conditionally to $X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_p$

to the prediction of Y

Obviously, if multicollinearity

contribution of $X_i \approx 0$.

Let see an example :

"The price of the cars."

(13)

How to detect the multicollinearity?

- Principal Component Analysis

Find linear combinations of X_i 's
of maximal variance.

$$c = u_1 X_1 + \dots + u_p X_p, \quad u = \begin{pmatrix} u_1 \\ \vdots \\ u_p \end{pmatrix} \in \mathbb{R}^p$$

such that

$$u = \arg \max_{\substack{u \in \mathbb{R}^p \\ \|u\| = 1}} V(X \cdot u) = V(c).$$

u is the eigen vector associated
to the largest eigen value of matrix
of variance-covariance of X 's

$$\underline{Vu = \lambda u}$$

14

Let $\{(\lambda_1, u_1), \dots, (\lambda_p, u_p)\}$ be the set of eigen values/vectors of the covariance matrix V , such that

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0.$$

Let also define the i^{th} principal component by

$$c_i = X^T \cdot u_i = X_1 \cdot u_{1,1} + \dots + X_p \cdot u_{1,p}$$

Then, it is well known that :

$$1) \lambda_1 + \lambda_2 + \dots + \lambda_p = \sum_{i=1}^p V(x_i) \quad \left(\begin{array}{l} = p \\ \uparrow \\ \text{scaled} \\ \text{data} \end{array} \right)$$

$$2) \begin{cases} V(c_i) = \lambda_i, & \forall i = 1, \dots, p \\ E(c_i) = 0 \end{cases}$$

$$3) X = c_1 \cdot u_1 + c_2 \cdot u_2 + \dots + c_p \cdot u_p$$

(expansion formula).

(15)

By 2), if there exists $k \in \{1, \dots, p\}$
such that

$$\lambda_k \approx 0$$

then

$$\forall (c_i) \approx 0 \Leftrightarrow c_i \approx 0$$

$$\Leftrightarrow$$

$$\underline{X_1 u_{k,1} + X_2 u_{k,2} + \dots + X_p u_{k,p} \approx 0}$$

multicollinearity.

See on the example. (car example).

(16)

Solutions to multicollinearity for linear model estimation

Remark that, because of the expansion formula 2), the linear space spanned by X_i 's is the same as that spanned by the principal components c_i 's.

Principal component Regression (PCR)

Regression of Y using the principal components with variance > 0 .

$$\hat{Y} = \gamma_1 c_1 + \gamma_2 c_2 + \dots + \gamma_q c_q$$

17

Notice that

$$V(\hat{\delta}_i) = \frac{\sigma^2}{n} \cdot \frac{1}{\lambda_i} \quad \text{and}$$

$$\hat{\beta}_j = \hat{\delta}_1 u_{1,j} + \hat{\delta}_2 u_{2,j} + \dots + \hat{\delta}_q u_{q,j}$$

$$V(\hat{\beta}_j) = \frac{\sigma^2}{n} \cdot \sum_{i=1}^q \frac{u_{ij}^2}{\lambda_i}$$

Remark that small values of λ_i \approx multicollinearity

\Downarrow
explosion of the variance of $\hat{\beta}_j$'s.

Idea: Keep only the principal comps. with large variances (λ_i 's).

But: the least squares criterion maximizes R^2 and the most explanatory c.p. are not necessarily correlated to Y !

①⑧ Since the c_i 's are uncorrelated

$$R^2(Y, \underbrace{\gamma_1 c_1 + \dots + \gamma_q c_q}_{\hat{Y}}) = \underbrace{R^2(Y, c_1)} + \dots + \underbrace{R^2(Y, c_q)}$$

• See the example on cars.

- use stepwise selection of p.c.'s

- use the first q components selected
by cross-validation procedure
(package `pls`, function `pcr`)

(19)

So, ideally, one wants "principal components" which:

- are highly correlated with Y (R^2)
- explain large amount of information from X . (λ)

This is not always the case and a compromise between the goodness of fit (R^2) and the stability of the regression coefficients ($V(\hat{\beta}_j)$) must be done!

It is a difficult task!

- The PLS regression builds such components maximizing the covariance instead the R^2 .

(20) Let look for a component

$$t = u_1 X_1 + u_2 X_2 + \dots + u_p X_p$$

Such that

$\text{Cov}^2(t, Y)$ is maximized Tucker
criterion

under the constraint $\sum_{i=1}^p u_i^2 = 1$

Remark that, because

$$\text{Cov}^2(t, Y) = \underline{R^2(t, Y)} \cdot \underline{V(t)} \cdot V(Y)$$

one maximizes simultaneously both

$$\underline{R^2(Y, t)}$$

and

$$\underline{V(t)}$$

↓
link with Y

↓
link with X

This component which maximizes a
such criterion (Tucker) is called
first PLS component

(21) Let denote by t_1 this first PLS component:

$$t_1 = X_1 u_{11} + X_2 u_{12} + \dots + X_p u_{1p}$$

The weights $\{u_{1i}\}_{i=1, \dots, p}$ which satisfy the Tucker criterion are given by:

$$u_{1,i} = \frac{E(Y \cdot X_i)}{\sqrt{\sum_{j=1}^p E^2(Y \cdot X_j)}}$$

Of course, t_1 "explains" X , in general, less than the first principal component c_1 , but is better for the prediction of Y :

$$\text{de Jong (1993): } \underline{R^2(Y, t_1) \geq R^2(Y, c_1)}$$

(22) We are looking now for a second PLS component, t_2 .

For this purpose we perform the following simple linear regressions:

Put $X^{(0)} = X$ and $Y^{(0)} = Y$

$$\begin{cases} X_i^{(0)} = \underline{p_{1,i}} t_1 + \underbrace{X_i^{(1)}}_{\text{residual}}, & i=1, \dots, p. \\ Y^{(0)} = \underline{c_1} t_1 + \underbrace{Y^{(1)}}_{\text{residual}} \end{cases}$$

t_2 is built in the same way as t_1 but using $X^{(1)}$ and $Y^{(1)}$
residuals after regression on t_1 .

$$\begin{cases} X^{(1)} = (X_1^{(1)}, X_2^{(1)}, \dots, X_p^{(1)}) \\ Y^{(1)} \end{cases}$$

(23)

$$t_2 = u_{21}X_1^{(1)} + u_{22}X_2^{(1)} + \dots + u_{2p}X_p^{(1)}$$

with $u_2 = (u_{21}, \dots, u_{2p})$ such that
the Tucker criterion for $X^{(1)}$ and $Y^{(1)}$

$$\max \text{Cor}^2(t_2, Y^{(1)}) = \max \text{Cor}^2\left(\sum_{i=1}^p u_{2i}X_i^{(1)}, Y^{(1)}\right)$$

$$u_{2i} = \frac{\mathbb{E}(Y^{(1)}X_i^{(1)})}{\sqrt{\sum_{j=1}^p \mathbb{E}^2(Y^{(1)}X_j^{(1)})}}, \quad i=1, \dots, p.$$

Notice that:

$$\begin{cases} \cdot X_i^{(1)} \perp t_1 & (\mathbb{E}(X_i^{(1)}t_1) = 0) \\ \cdot Y^{(1)} \perp t_1 & (\mathbb{E}(Y^{(1)}t_1) = 0) \end{cases}$$

$$\rightarrow t_1 \perp t_2 \quad (\mathbb{E}(t_1, t_2) = 0).$$

(24)

Remark also that

- t_1 is linear combination of X_i 's
- t_2 is linear combination of $X_i^{(1)}$'s but also of X_i 's

$$t_2 = \sum_{i=1}^P u_{2i} X_i^{(1)} = \underbrace{\sum_{i=1}^P u_{2i} (X_i - p_{2,i} t_1)}_{\text{linear comb of } X_i \text{'s}}$$

In the similar way as for t_1 , we compute new residuals using t_2 :

$$\begin{cases} X_i^{(1)} = \underbrace{p_{2,i} t_2}_{\text{residual}} + \underbrace{X_i^{(2)}}_{\text{residual}} \\ Y^{(2)} = \underbrace{r_2 t_2}_{\text{residual}} + \underbrace{Y^{(2)}}_{\text{residual}} \end{cases}$$

(25)

The computation procedure of t_h and of $X^{(h)}$ and $Y^{(h)}$ is called the h^{th} step of the PLS regression of Y on X . (In our example $h=2$).

Expansion formulas and linear approximation

$h=2$:

$$\begin{cases} X_i = p_{1,i} t_1 + X_i^{(1)} \\ Y = \kappa_1 t_1 + Y^{(1)} \end{cases} \quad \text{and} \quad \begin{cases} X_i^{(1)} = p_{2,i} t_2 + X_i^{(2)} \\ Y^{(1)} = \kappa_2 t_2 + Y^{(2)} \end{cases}$$

$$\begin{cases} X_i = \underbrace{p_{1,i} t_1 + p_{2,i} t_2}_{\text{prediction } f(x)} + \underbrace{X_i^{(2)}}_{\text{residual}} \\ Y = \underbrace{\kappa_1 t_1 + \kappa_2 t_2}_{\text{prediction } f(x)} + \underbrace{Y^{(2)}}_{\text{residual}} \end{cases}$$

(26)

PLS regression is an iterative method.

At step \underline{h} we have:

$\{t_1, t_2, \dots, t_h\}$: the PLS components
(the first h)

⚠ $\forall i \neq j$ t_i and t_j are uncorrelated!

We have also the following expansion formulas:

$$\begin{cases} X_i = p_{1,i} t_1 + p_{2,i} t_2 + \dots + p_{h,i} t_h + X_i^{(h+1)} \\ Y = \underbrace{c_1 t_1 + c_2 t_2 + \dots + c_h t_h}_{\text{PLS prediction at step } h} + \underbrace{Y^{(h+1)}}_{\text{residual.}} \end{cases}$$

A question: how large is h ?

- of course, because $\{t_i\}_{i=1, \dots, h}$ are uncorrelated

$$\underline{h \leq \dim(L(x))}$$

(27)

More precisely, if $L_Y(x)$ is the smallest subspace of $L(X)$ which contains $\hat{Y} = P_x(Y)$, then

$$\underline{h \leq \dim(L_Y(x))}$$

In practice, h is chosen by cross-validation

Root Mean Squared Error of Prediction:
(RMSEP)

$$\text{RMSEP} = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$$

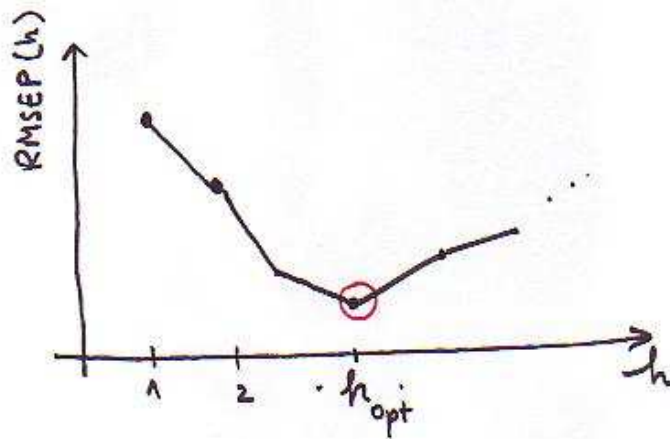
In general \hat{Y}_i is computed with the model built without the observation "i"
(leave-one-out cross-validation).

For fixed h , one has

$$\text{RMSEP}(h).$$

(28)

One can choose h which minimizes the function $RMSEP(h)$



Tenenhaus proposed to stop the iterative process when the $(h+1)$ PLS component does not contribute significantly to the prediction:

- if $\alpha \in (0, 1)$ - typically 0.9 or 0.95, the process stops at step h if

$$RMSEP(h+1) \geq \alpha \cdot RMSEP(h)$$

\Leftrightarrow

$$\frac{RMSEP(h) - RMSEP(h+1)}{RMSEP(h)} \leq \underline{1-\alpha}$$

(the relative contribution is less than $1-\alpha$).