collinearity

big K. small n. (X'X) (K+1)*(K+1) ⇒ B = <XXXXX is not unique

Multicollinearity in regression refers to the case when one predictor variable in a multiple regression model can be linearly predicted from the others with a substantial degree of accuracy.

> 把其中一些 Variable 五样 (Covariate) 美联性较大的

Detection

在什么情况下,不能有失线性

- 1. Significant correlations between pairs of independent variables in the model
- 2. Nonsignificant t-tests for all (or nearly all) the individual eta parameters when the F-test for overall model adequacy $H_0\colon eta_1=eta_2=\dots=eta_k=0$ is significant
- 3. Opposite signs (from what is expected) in the estimated parameters
- 4. A variance inflation factor (VIF) for a eta parameter greater than 10, where

Criterion: $R_i^2 > 0.9 \Leftrightarrow \text{CVLF} > 10^{(VIF)_i} = \frac{1}{1 - R_i^2}$ i = 1, 2, ..., k

and R_i^2 is the multiple coefficient of determination for the model

 $E(x_i) = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_{i-1} x_{i-1} + \alpha_{i+1} x_{i+1} + \dots + \alpha_k x_k \qquad \begin{cases} 2 & \\ 1 &$

Condition number

 $\kappa = \sqrt{\frac{\lambda_1}{\lambda_p}}$ some eigenvalue = 0 if vank(X) < kH $\Rightarrow k \Rightarrow \infty$

 $\sqrt{\frac{\lambda_1}{\lambda_i}} \quad \text{where } i = 2, ..., p \qquad \sqrt{\frac{\lambda_1}{\lambda_i}} \rightarrow \infty \quad \text{for some } \underline{i}$

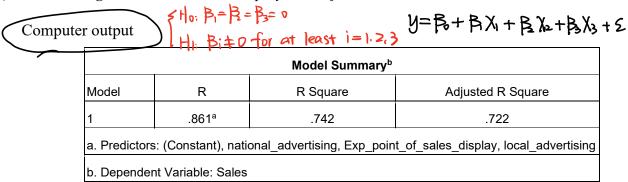
Multicollinear problem: Condition number or indexes \geq 30

(note in some references, condition index is also called condition number)

Motivation example

An assistant in the district sales office of a national cosmetics firm obtained data on advertising expenditures and sales last year in the district's 44 territories. X1 denotes expenditures for point-of-sales displays in beauty salons and department stores (in thousand dollars), and X2 and X3 represent the corresponding expenditures for local media advertising and prorated share of national media advertising, respectively. Y denotes sales (in thousand cases). The assistant was instructed to estimate the increase in expected sales when X1 is increased by 1 thousand dollars and X2 and X3 are held constant, and was told to use an ordinary multiple regression model with linear terms for the predictor variables.

a) State the regression model to be employed and fit it to the data.



	ANOVA ^b F 70 F(3,47)									
М	Model Sum			of Squares	df	Mean Square		F	Sig.	
1		Regression	38	32.659	3	127.5	53	38.279	.000ª)
		Residual	13	33.286	40	3.332	2		/	/
		Total	5	15.945	43				Pralu	e => reject to Ho
			T			Pvalu	e → Ho; B;	j= 0 V.S. H : B	j‡0	
				coefficients	Std. Error	t	sig	\	/IF	
1	(Cor	nstant)		1.023	1.203	.851	.400			
	Exp_	_point_of_sales_d	display	.966	.709	1.362	.181	20	0.072 X~K	+1/3 -> Ki=> VIF
	loca	_advertising		.629	.778	.808	.424	20).716 X ₂ ~X	+X3 -> R12 -> VIF +X3 -> R22 -> VIF
	national_advertising		.676	.356	1.900	.065	1.	.218		

Refer to the above output, the regression equation is.

$$\hat{Y} = 1.023 + .966X_1 + .629X_2 + .676X_3$$

b) Test whether there is a regression relation between sales and the three predictor variables (use $\alpha = 0.05$). State the alternatives, decision rule, and conclusion.

$$E(y) = \beta_0 + \beta_1 \ X_1 + \beta_2 \ X_2 + \beta_3 \ X_3 \qquad \qquad H_0: \ \beta_1 = \beta_2 = \beta_3 = 0 \qquad H_a: \ not \ all \ \beta_k = 0 \ (k = 1, \ 2, \ 3).$$
 Test Statistics: $F = 38.279$ and $p\text{-value} < 0.001$.

Therefore, conclude H_a. The model is useful for the prediction of sales.

c) Test for each of the regression coefficients (equal to zero?) individually (use $\alpha = 0.05$ each time). Do the conclusions of these tests correspond to that obtained in part (b)?

Refer to part (a) output, all regression coefficients are not significant at the given α level. Therefore, these tests do not yield the same conclusion as in part (b). This is a consequence of the multi-collinearity problem.

d) Obtain the correlation matrix of the X variables and comment on the suitability of the data for the research objective.

From the correlation matrix below, we observe that the independent variables (X1 and X2) are highly correlated and the regression model is therefore not quite appropriate.

Correlations (ΜΑΤΥΊΧ)								
			Exp point-of-	Local	National			
		Sales	Sales display	advertising	advertising			
Sales	Pearson Correlation	1	.842**	.842**	.474**			
	Sig. (2-tailed)		.000	.000	.001			
Exp point-of-sales display	Pearson Correlation	.842**	1	.974**	.376*			
	Sig. (2-tailed)	.000		.000	.012			
Local advertising	Pearson Correlation	.842**	.974**	1	.410**			
	Sig. (2-tailed)	.000	.000		.006			
National advertising	Pearson Correlation	.474**	.376*	.410**	1			
	Sig. (2-tailed)	.001	.012	.006				
**. Correlation is significant at the 0.01 level (2-tailed), *. Correlation is significant at the 0.05 level (2-tailed).								

e) Obtain the three variance inflation factors. What do these suggest about the effects of multicollinearity here?

$$(V \text{ IF})_1 = 20.072$$

 $(V \text{ IF})_2 = 20.716$
 $(V \text{ IF})_3 = 1.218$

The problem is quite serious since two of the VIF are much larger than 10.

f) The assistant eventually decided to drop variables X2 from the model to clear up the picture. Fit the assistant's revised model. Is the assistant now in a better position to achieve the research objective?

ANOVAb									
Model		Sum of Squares	df	Mean Square	F	Sig.			
2	Regression	380.481	2	190.241	57.579	.000ª			
	Residual	135.464	41	3.304					
	Total	515.945	43						

a. Predictors: (Constant), national_advertising, Exp_point_of_sales_display

b. Dependent Variable: Sales

		Coefficient	Std. Error	t	sig	VIF
2	(Constant)	1.017	1.198	.849	.401	
XI	Exp_point_of_sales_display	1.522	.170	8.948	.000	1.165
X2	national_advertising	.736	.346	2.125	.040	1.165

Refer to the above output, the regression equation is

$$\hat{\mathbf{Y}} = 1.017 + 1.522X_1 + .736X_3$$

The VIF indicates that the problem of multicollinearity disappears (much less than 10).

Consequences of multicollinearity

- 1. Difficult to test individual regression coefficients due to inflated standard errors.
- 2. Unstable coefficient estimates, sensitive to small change in the model.

Remedial measures

- 1. Drop one or several highly correlated independent variables (or by stepwise regression to select appropriate variables)
- 2. Combine variables (dimension reduction by the use of methods such as principle component procedures)
- 3. Shrinkage methods such as ridge regression

$$\hat{\beta} = (\chi \chi)^{\dagger} \chi' \chi$$

$$\Rightarrow \hat{\beta} = (\chi \chi + \delta I)^{\dagger} \chi' \chi \qquad \underline{5 > 0}$$

$$\sim \text{ridge}$$

Another Example (extracted from Linear Models with R, Julian J. Faraway)

Car drivers like to adjust the sear position for their own comfort. Car designers would find it helpful to know where different drivers will position the seat depending on their size and age.

```
Sample size: 38 (drivers)
Dependent variable: hipcenter (the horizontal distance of the midpoint of the hips
                   from a fixed location in the car in mm)
Independent variables:
           Age
           Weight
           HtShoes (height with shoes)
           Ht (height without shoes)
           Seated (seated height)
           Arm (arm length)
           Thight (thigh length)
           Leg (lower leg length)
> library(faraway)
> data(seatpos)
> names(seatpos)
[1] "Age"
                "Weight"
                            "HtShoes"
                                       "Ht"
                                                   "Seated"
                                                               "Arm"
                                                                           "Thigh"
                "hipcenter"
[8] "Leg"
> reg1 <- lm(hipcenter~., seatpos)</pre>
> summary(reg1)
lm(formula = hipcenter ~ ., data = seatpos)
Residuals:
    Min
             10 Median
                            30
                                   Max
-73.827 -22.833 -3.678 25.017 62.337
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 436.43213 166.57162 2.620
                                          0.0138 *
                                  1.360
Age
             0.77572
                        0.57033
                                          0.1843
Weight
             0.02631
                        0.33097
                                  0.080
                                         0.9372
HtShoes
            -2.69241
                       9.75304 -0.276
                                          0.7845
Ht
             0.60134 10.12987
                                 0.059
                                          0.9531
                                 0.142
                                          0.8882
Seated
             0.53375
                        3.76189
                        3.90020 -0.341
                                         0.7359
Arm
            -1.32807
Thigh
            -1.14312
                        2.66002
                                 -0.430 0.6706
Leg
            -6.43905
                        4.71386 -1.366
                                         0.1824
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' ' 1
Residual standard error: 37.72 on 29 degrees of freedom
Multiple R-squared: 0.6866, Adjusted R-squared: 0.6001
F-statistic: 7.94 on 8 and 29 DF, p-value: 1.306e-05
```

*****F-test significant, but NOT individual t-test.****

⁵

```
>round(cor(seatpos),4)
             Age Weight HtShoes
                                    Ht Seated
                                                   Arm
                                                        Thigh
                                                                  Leg hipcenter
Age
          1.0000 0.0807 -0.0793 -0.0901 -0.1702 0.3595
                                                       0.0913 -0.0423
                                                                        0.2052
Weight
          0.0807 1.0000 0.8282 0.8285 0.7756 0.6976 0.5726 0.7843
                                                                        -0.6403
HtShoes
         -0.0793 0.8282 1.0000 0.9981 0.9297 0.7520 0.7249 0.9084
                                                                        -0.7966
         -0.0901 0.8285 0.9981 1.0000 0.9282 0.7521 0.7350 0.9098
                                                                        -0.7989
Ηt
Seated
         -0.1702
                 0.7756 0.9297 0.9282 1.0000 0.6252
                                                       0.6071
                                                               0.8119
                                                                        -0.7313
                  0.6976 0.7520 0.7521 0.6252
Arm
          0.3595
                                                1.0000
                                                       0.6711
                                                               0.7538
                                                                        -0.5851
Thigh
          0.0913 0.5726 0.7249 0.7350 0.6071 0.6711
                                                       1.0000 0.6495
                                                                        -0.5912
         -0.0423 0.7843 0.9084 0.9098 0.8119 0.7538 0.6495
                                                                        -0.7872
Leg
                                                               1.0000
hipcenter 0.2052 -0.6403 -0.7966 -0.7989 -0.7313 -0.5851 -0.5912 -0.7872
                                                                        1.0000
```

Some very large pairwise correlations

```
> x <- model.matrix(reg1) [,-1] > YEMOVE

> e <- eigen(t(x) %*% x)

> e$val

[1] 3.653671e+06 2.147948e+04 9.043225e+03 2.989526e+02 1.483948e+02 8.117397e+01

5.336194e+01

[8] 7.298209e+00

> sqrt(e$val[1]/e$val) # compute condition indexes

[1] 1.00000 13.04226 20.10032 110.55123 156.91171 212.15650 261.66698 707.54911
```

Very large condition numbers

VIFs of HtShoes and Ht are extremely high

Q: how to use Im to calculate Riz and LVIF);

Rerun the regression with Age, Weight and Ht only

```
> reg2 <- lm(hipcenter~Age+Weight+Ht,seatpos)</pre>
> x <- model.matrix(reg2) [,-1]</pre>
> e <- eigen(t(x) %*% x)
> sqrt(e$val[1]/e$val) # compute condition indexes
[1] 1.00000 11.50837 15.50904
> vif(x)
     Age
           Weight
1.093018 3.457681 3.463303
> summary(reg2)
Call:
lm(formula = hipcenter ~ Age + Weight + Ht, data = seatpos)
Residuals:
    Min
             1Q Median
                                    Max
                             3Q
-91.526 -23.005
                  2.164 24.950 53.982
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 528.297729 135.312947
                                    3.904 0.000426 ***
              0.519504
                         0.408039
                                    1.273 0.211593
Weight
              0.004271
                         0.311720
                                    0.014 0.989149
                         0.999056 -4.216 0.000174 ***
Ht
             -4.211905
Signif. codes: 0 '***, 0.001 '**, 0.01 '*, 0.05 '., 0.1 ', 1
Residual standard error: 36.49 on 34 degrees of freedom
Multiple R-squared: 0.6562, Adjusted R-squared: 0.6258
F-statistic: 21.63 on 3 and 34 DF, p-value: 5.125e-08
> round(cor(x),3)
          Age Weight
                         Ht
        1.000 0.081 -0.090
Age
Weight 0.081 1.000 0.829
Ht
       -0.090 0.829 1.000
```