

Exam Practice Questions

Calculate the estimators I:

Consider the following dataset, where Y is the dependent variable and X_1 and X_2 are the independent variables:

Given Data:			
Observation	Y	X_1	X_2
1	5	1	2
2	10	2	3
3	15	3	4
4	20	4	5

We want to fit a linear regression model of the form $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$

Compute the least squares estimates of β using the matrix formula:

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

Calculate the $X^T X$ and $X^T Y$

Result

Now, compute $X^T X$:

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 5 \end{bmatrix} \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 3 \\ 1 & 3 & 4 \\ 1 & 4 & 5 \end{bmatrix} = \begin{bmatrix} 4 & 10 & 14 \\ 10 & 30 & 40 \\ 14 & 40 & 54 \end{bmatrix}$$

Next, compute $X^T Y$:

$$X^T Y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 5 \end{bmatrix} \begin{bmatrix} 5 \\ 10 \\ 15 \\ 20 \end{bmatrix} = \begin{bmatrix} 50 \\ 150 \\ 200 \end{bmatrix}$$

We got the β_0 , β_1 and β_2 when we multiply $(X^T X)^{-1}$ and $X^T Y$.

Calculate the estimators II:

The inverse of $X^T X$ is given:

The inverse of $X^T X$ is:

$$(X^T X)^{-1} = \begin{bmatrix} 1.5 & -0.5 & 0 \\ -0.5 & 0.5 & 0 \\ 0 & 0 & 0.5 \end{bmatrix}$$

Find the estimators

Result

Now, multiply $(X^T X)^{-1}$ by $X^T Y$:

$$\beta = \begin{bmatrix} 1.5 & -0.5 & 0 \\ -0.5 & 0.5 & 0 \\ 0 & 0 & 0.5 \end{bmatrix} \begin{bmatrix} 50 \\ 150 \\ 200 \end{bmatrix} = \begin{bmatrix} 0 \\ 5 \\ 10 \end{bmatrix}$$

The estimated coefficients are:

$$\beta_0 = 0, \quad \beta_1 = 5, \quad \beta_2 = 10$$

Thus, the estimated regression equation is:

$$Y = 0 + 5X_1 + 10X_2$$

Given our estimated model, calculate the MS reg and MS res:

Result

$$SS_{reg} = \sum (\hat{Y}_i - \bar{Y})^2$$

Compute SS_{reg} :

$$SS_{reg} = (25 - 12.5)^2 + (40 - 12.5)^2 + (55 - 12.5)^2 + (70 - 12.5)^2$$

$$SS_{reg} = (12.5)^2 + (27.5)^2 + (42.5)^2 + (57.5)^2$$

$$SS_{reg} = 156.25 + 756.25 + 1806.25 + 3306.25 = 6025$$

Residual Sum of Squares (SS_{res}):

$$SS_{res} = \sum (Y_i - \hat{Y}_i)^2$$

Compute SS_{res} :

$$SS_{res} = (5 - 25)^2 + (10 - 40)^2 + (15 - 55)^2 + (20 - 70)^2$$

$$SS_{res} = (-20)^2 + (-30)^2 + (-40)^2 + (-50)^2$$

$$SS_{res} = 400 + 900 + 1600 + 2500 = 5400$$

- **Degrees of Freedom for Regression (df_{reg}):** Number of predictors = 2 (X_1 and X_2).
- **Degrees of Freedom for Residual (df_{res}):** $n - p - 1 = 4 - 2 - 1 = 1$.

Mean Square Regression (MS_{reg}):

$$MS_{reg} = \frac{SS_{reg}}{df_{reg}} = \frac{6025}{2} = 3012.5$$

Mean Square Residual (MS_{res}):

$$MS_{res} = \frac{SS_{res}}{df_{res}} = \frac{5400}{1} = 5400$$

Now form the hypothesis for our model, is our model significant? (F test):

Result

Step 1: Calculate the F-statistic (F_{obs})

The F-statistic is calculated as:

$$F_{obs} = \frac{MS_{reg}}{MS_{res}}$$

From the previous calculations:

- $MS_{reg} = 3012.5$
- $MS_{res} = 5400$

Now, compute F_{obs} :

$$F_{obs} = \frac{3012.5}{5400} \approx 0.5579$$

Step 2: Determine the Critical Value from the F-distribution Table

The critical value for the F-distribution depends on:

- **Degrees of Freedom for Regression (df_{reg}): 2**
- **Degrees of Freedom for Residual (df_{res}): 1**
- **Confidence Level (α): 0.05**

Using an F-distribution table or calculator, the critical value $F_{(2,1,0.05)}$ is approximately **18.51**.

Step 3: Compare F_{obs} to the Critical Value

- $F_{obs} \approx 0.5579$
- Critical Value $F_{(2,1,0.05)} = 18.51$

Since $F_{obs} = 0.5579 < 18.51$, we **fail to reject the null hypothesis** at the 0.05 confidence level.

Step 4: Interpretation

- The null hypothesis H_0 states that all regression coefficients (β_1 and β_2) are zero, meaning the predictors X_1 and X_2 do not significantly explain the variation in Y .
 - Since F_{obs} is less than the critical value, we do not have sufficient evidence to reject H_0 . This suggests that the regression model does not provide a significant fit to the data at the 0.05 confidence level.
-

Final Answer:

- **F-statistic (F_{obs}):** 0.5579
 - **Critical Value ($F_{(2,1,0.05)}$):** 18.51
 - **Conclusion:** Fail to reject the null hypothesis H_0 at the 0.05 confidence level. The regression model is not statistically significant.
-

Does adding X3, X4, X5, X6, and X7 to the model have a significant impact on improving the model's predictions?

STAT 5004

Example 11.1

Page 6

MLR EXAMPLE 11.1 - GAS MILEAGE

4

Model: MODEL2
Dependent Variable: MPG

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F
Model	2	1339.21617	669.60809	66.903	0.0001
Error	29	290.25101	10.00866		
C Total	31	1629.46719			
Root MSE		3.16365	R-square	0.8219	
Dep Mean		20.29062	Adj R-sq	0.8096	
C.V.		15.59166			

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T
INTERCEP	1	33.491421	1.47544969	22.699	0.0001
SIZE	1	-0.040148	0.00990172	-4.055	0.0003
HP	1	-0.016611	0.02250347	-0.738	0.4663

Model: MODEL1
Dependent Variable: MPG

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F
Model	7	1466.87390	209.55341	30.932	0.0001
Error	24	162.59329	6.77472		
C Total	31	1629.46719			
Root MSE		2.60283	R-square	0.9002	
Dep Mean		20.29062	Adj R-sq	0.8711	
C.V.		12.82774			

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T
INTERCEP	1	45.026054	12.37472908	3.639	0.0013
SIZE	1	0.077614	0.03039037	2.554	0.0174
HP	1	-0.032959	0.01916261	-1.720	0.0983
WEIGHT	1	-0.008264	0.00228236	-3.621	0.0014
SHAPE	1	1.283181	2.77643544	0.462	0.6481
CYLIN	1	-3.169670	1.44712974	-2.190	0.0385
TRANS	1	0.510739	3.58673059	0.142	0.8880
SPEEDS	1	1.869099	2.78918048	0.670	0.5092

Variable	DF	Type III SS
INTERCEP	1	89.690772
SIZE	1	44.187532
HP	1	20.040986
WEIGHT	1	88.815898
SHAPE	1	1.447078
CYLIN	1	32.501586
TRANS	1	0.137370
SPEEDS	1	3.042306

Result

$$R(\beta_3, \beta_4, \beta_5, \beta_6, \beta_7 \mid \beta_1, \beta_2) = \text{SSreg full} - \text{SSreg sub}$$

$$\text{Fobs} = (\text{SSreg full} - \text{SSreg sub} / p) / \text{MSres full model}$$

To form this hypothesis we first look at the SSreg of the sub-model = 1339,
Next, we look at the SSreg of the full model containing X3, X4, X5, X6 = 1466

1. SSreg full model - SSreg sub model = 1466 - 1339 = 127
2. Divide it by the number of parameters in the model = 127 / 5 = 25.4
3. Divide it by MSres full = 25.4 / 6.775 = 3.75 = Fobs
4. Look at the F table for (p, n-p-1) = F(5,24) = 2.62
5. Fobs (3.75) > Ftable (2.62), Reject the null hypothesis
6. Rejecting the H0 means adding these predictors will improve the model's prediction quality significantly.

$$\begin{aligned}
 \text{F test: } R(\beta_3, \beta_4, \beta_5, \beta_6, \beta_7 | \beta_1, \beta_2) &= R(\beta_1, \beta_2, \dots, \beta_7) - R(\beta_1, \beta_2) \\
 &= 1466.874 - 1339.216 \\
 &= 127.658 \\
 F_{\text{obs}} &= \frac{R(\beta_3, \beta_4, \beta_5, \beta_6, \beta_7 | \beta_1, \beta_2) / 5}{MS_{\text{res, full}}} = \frac{127.658 / 5}{6.775} \\
 &= 3.769 \sim F(5, 24) \\
 F(5, 24)_{\text{crit}} &= 2.62 \Rightarrow \text{reject } H_0 \quad (0.01 < p < 0.05)
 \end{aligned}$$

Can we remove β_4 ? β_3 ? β_2 ? (MSres = 0.0499, n = 17)

d) Type I SS = sequential SS

$$\begin{aligned}
 R(\beta_1) &= R(\beta_1) &= 22.5590 \\
 R(\beta_2 | \beta_1) &= R(\beta_1, \beta_2) - R(\beta_1) &= 0.2825 \\
 R(\beta_3 | \beta_1, \beta_2) &= R(\beta_1, \beta_2, \beta_3) - R(\beta_1, \beta_2) &= 0.1353 \\
 R(\beta_4 | \beta_1, \beta_2, \beta_3) &= R(\beta_1, \beta_2, \beta_3, \beta_4) - R(\beta_1, \beta_2, \beta_3) &= 0.0041
 \end{aligned}$$

notes:

(i) adding these sequential SS yields:

$$R(\beta_1, \beta_2, \beta_3, \beta_4) = 22.9809 = SS_{\text{reg, full}}$$

Result

X4: Fobs = 0.0041 / 0.0499 = 0.082, F(1, 15) = 4.54, **Accept the H0** and **drop X4**,

X3: $F_{obs} = 0.1353 / 0.0499 = 2.710$, $F(1,15) = 4.54$, **Accept the H0** and **drop X3**,

X2: $F_{obs} = 0.2825 / 0.0499 = 5.65$, $F(1,15) = 4.54$, **Reject the H0** and **keep X2**.

e) example

Do we need X^4 ?

$$F_{obs} = \frac{R(\beta_4 | \beta_1, \beta_2, \beta_3) / 1}{MS_{res, full}} = \frac{.0041}{.0499} = .082 \sim F_{(1,15)}$$

$F_{(1,15), .95} = 4.54 \Rightarrow$ non-significant \Rightarrow drop X^4

Do we need X^3 ?

$$F_{obs} = \frac{R(\beta_3 | \beta_1, \beta_2) / 1}{MS_{res, full}} = \frac{.1353}{.0499} = 2.710 < 4.54$$

\Rightarrow non-significant \Rightarrow drop X^3

Do we need X^2 ?

$$F_{obs} = \frac{R(\beta_2 | \beta_1) / 1}{MS_{res, full}} = \frac{.2825}{.0499} = 5.657 > 4.54$$

\Rightarrow significant

\Rightarrow stop and include X^2 and X

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2$$

After deciding that 2nd order model is adequate, we must run the 2nd order model to obtain the final prediction equation.

input: $\left. \begin{array}{l} \text{PROC REG;} \\ \text{MODEL Y = X X2;} \end{array} \right\}$ to run 2nd order model

output from this is on p.3

Final prediction equation is:

$$\hat{Y} = .1776 + 2.0469X - .1756X^2$$

Can we remove the β_4 from the model?

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F
Model	7	1466.87390	209.55341	30.932	0.0001
Error	24	162.50220	6.77472		
C Total	31	1629.46719			
Root MSE	2.60283	R-square	0.9002		
Dep Mean	20.29052	Adj R-sq	0.8711		
C.V.	12.82774				

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T
INTERCEP	1	45.026054	12.37472908	3.639	0.0013
SIZE	1	0.077614	0.03039037	2.554	0.0174
HP	1	-0.032959	0.01916261	-1.720	0.0983
WEIGHT	1	-0.008264	0.00228236	-3.621	0.0014
SHAPE	1	-3.169670	1.44712974	-2.190	0.0385
CYLIN	1	0.510739	3.58673059	0.142	0.8880
TRANS	1	1.869099	2.78918048	0.670	0.5092
SPEEDS	1				

Variable	DF	Type II SS
INTERCEP	1	89.690772
SIZE	1	44.187532
HP	1	20.040986
WEIGHT	1	88.815898
SHAPE	1	1.447078
CYLIN	1	32.501586
TRANS	1	0.137370
SPEEDS	1	3.042306

Result

To see if we can remove the β_4 from the model we have to look at the predictors.

β_1 = SIZE, β_2 = HP, β_3 = WEIGHT, β_4 = SHAPE, ...

SHAPE predictor has 1.4471 Type I SS, which is equivalent to:

$R(\beta_4 \mid \beta_1, \beta_2, \beta_3, \beta_5, \beta_6, \beta_7) = SS \text{ full} - SS \text{ sub}$

F-obs = $(R(\beta_4 \mid \beta_1, \beta_2, \beta_3, \beta_5, \beta_6, \beta_7) / p) / MS \text{ res full}$ = $((SS \text{ full} - SS \text{ sub}) / p) / MS \text{ res full}$ model

$(1.4471 / 1) / 6.774 = 0.214$

Now we want to compare the F-obs with F-table = $F(1,24) = 4.26$

We now can compare $F_{\text{obs}} < F_{\text{table}}$, **Reject H_0** , We can drop β_4 because its contribution to the model is not significant.

Note: You can directly look at the "**T for H_0 : Parameter = 0**" value from the table to see if the parameter is significant according to Type II SS (partial SS)

!!! $F_{\text{obs}} = T_{\text{obs}}^2 = 0.462^2 = 0.214$



c) Type II SS = partial SS

Test of each vbl. as if it were the last to be added to the model. For each vbl. separately we ask whether it would be worth adding if all the other predictors were already in the model.

example 11.1 (p. 2 of output):

$$R(\beta_1 | \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7) = 44.1875$$

$$R(\beta_2 | \beta_1, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7) = 20.0410$$

(ii) each of these can be tested as a partial F test.

example: $H_0: \beta_4 = 0$ in full model

$$F_{\text{obs}} = \frac{R(\beta_4 | \beta_1, \beta_2, \beta_3, \beta_5, \beta_6, \beta_7) / 1}{MS_{\text{res, full}}} = \frac{1.4471 / 1}{6.7747}$$

$$= 0.214 \sim F_{(1, 24)}$$

$$F_{(1, 24), 0.95} = 4.26 \Rightarrow \text{accept } H_0 \text{ for } \alpha = 0.05$$

\Rightarrow we could delete X_4 from the full model without significant loss in predictability

(iii) SAS output contains

$$t_{\text{obs}} = \sqrt{F_{\text{obs}}} \text{ under heading}$$

("T FOR H_0 : PARAMETER = 0"). Note that for

$$X_4 (\text{SHAPE}), t_{\text{obs}} = .462 = \sqrt{0.214} = \sqrt{F_{\text{obs}}}$$

p-value is .6481, consistent with our conclusion

Explain

Explain the following SAS inputs

1. Proc PRINT;
2. Proc REG;
 - Model Y = X X2 X3 X4 / SS1;
3. Proc REG;
 - Model Y = X X2;
4. Proc REG;
 - MODEL MPG = SIZE WEIGHT SHAPE / P SS2;
 - RUN;

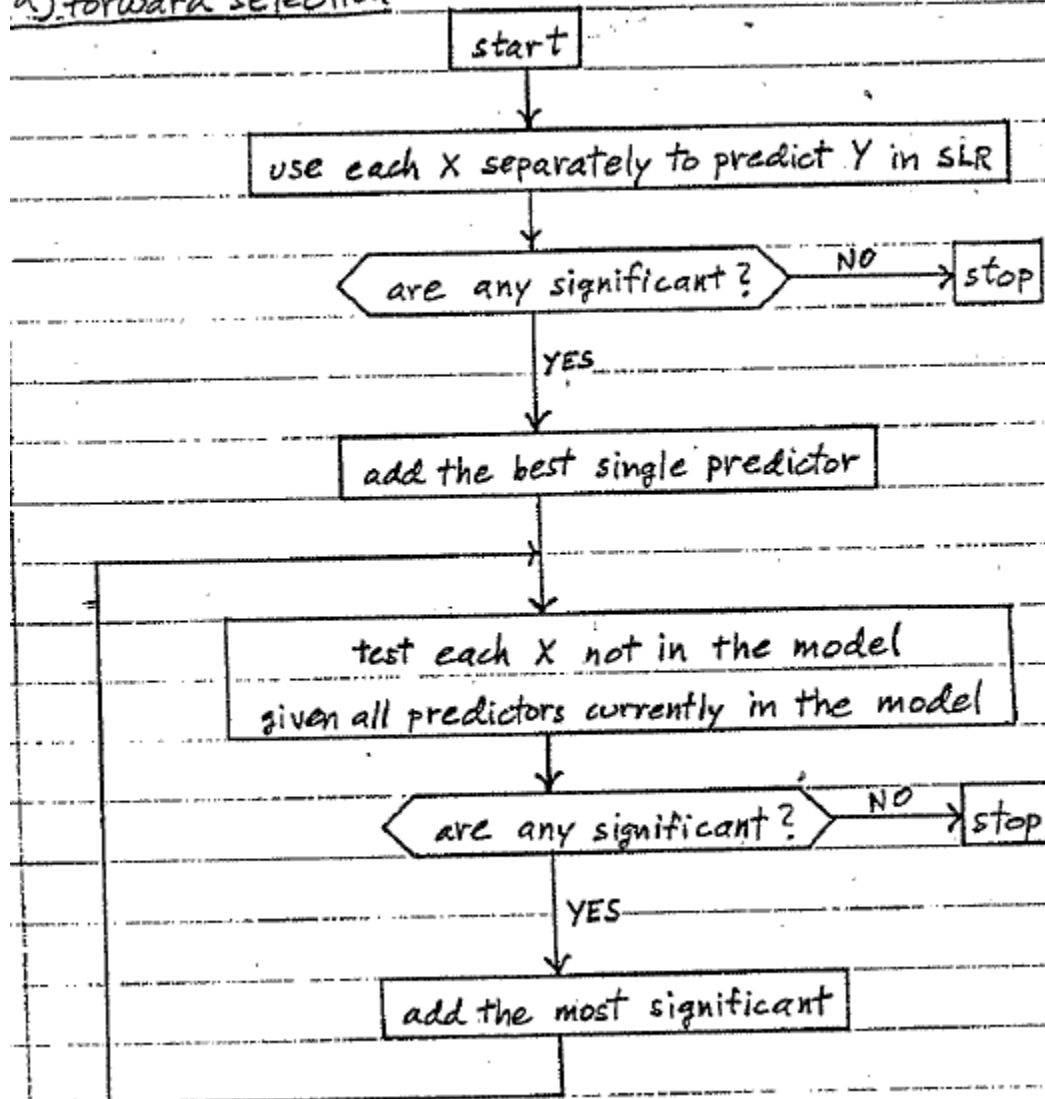
Result

1. Prints the data, the output contains the values of Y, X, X2, X3, X4 (ie: $X2 = X^2$)
 2. To examine 4th order model (X, X^2 , X^3 , X^4)
 3. To examine 2th order model (X, X^2)
 4. Builds a multiple linear regression model to predict Y (MPG) with X1 (Size), X2 (Weight) and X3 (Shape) features. Options: P = Requests predicted values (\hat{Y}) and residuals (e) for each observation. SS2: Provides the Type II Sum of Squares for the model.
-

Draw the flow chart for the Forward Selection:

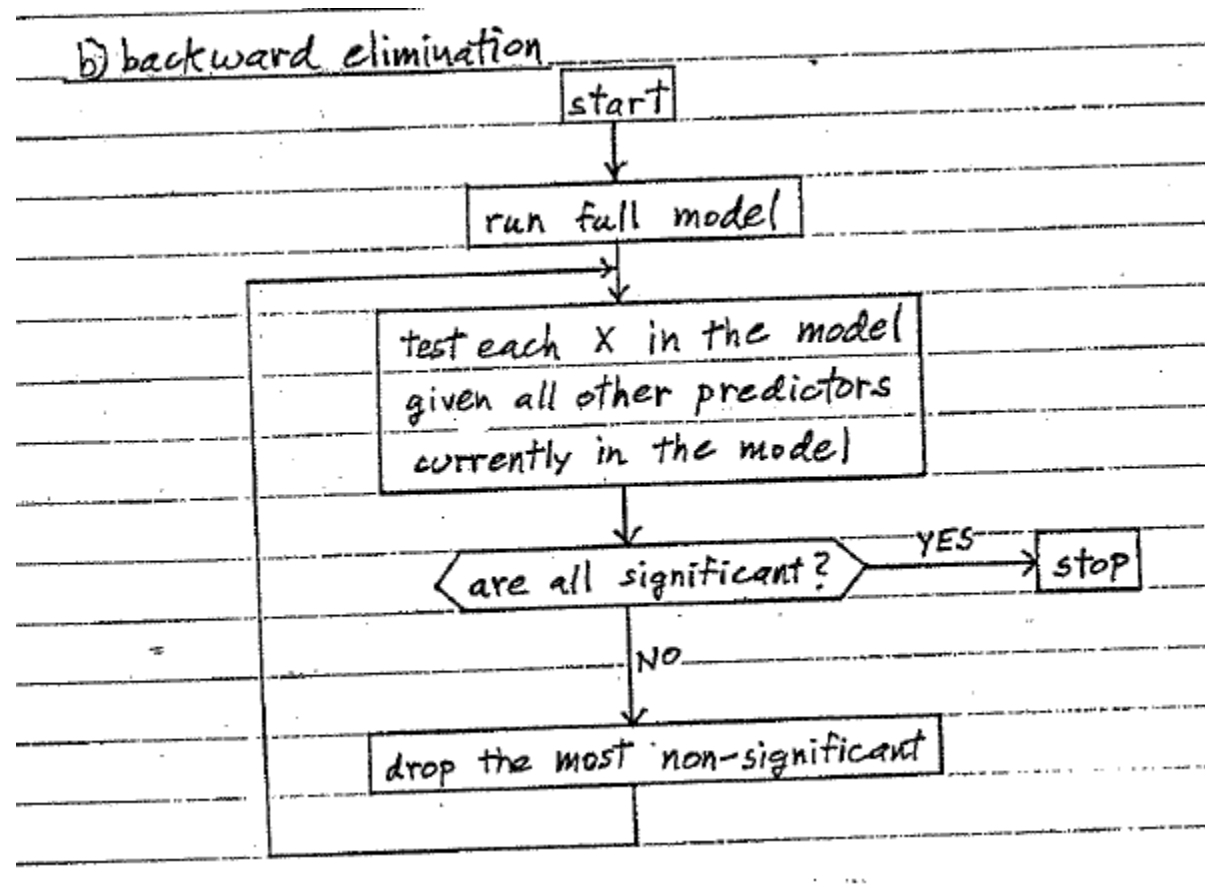
Result

a) forward selection



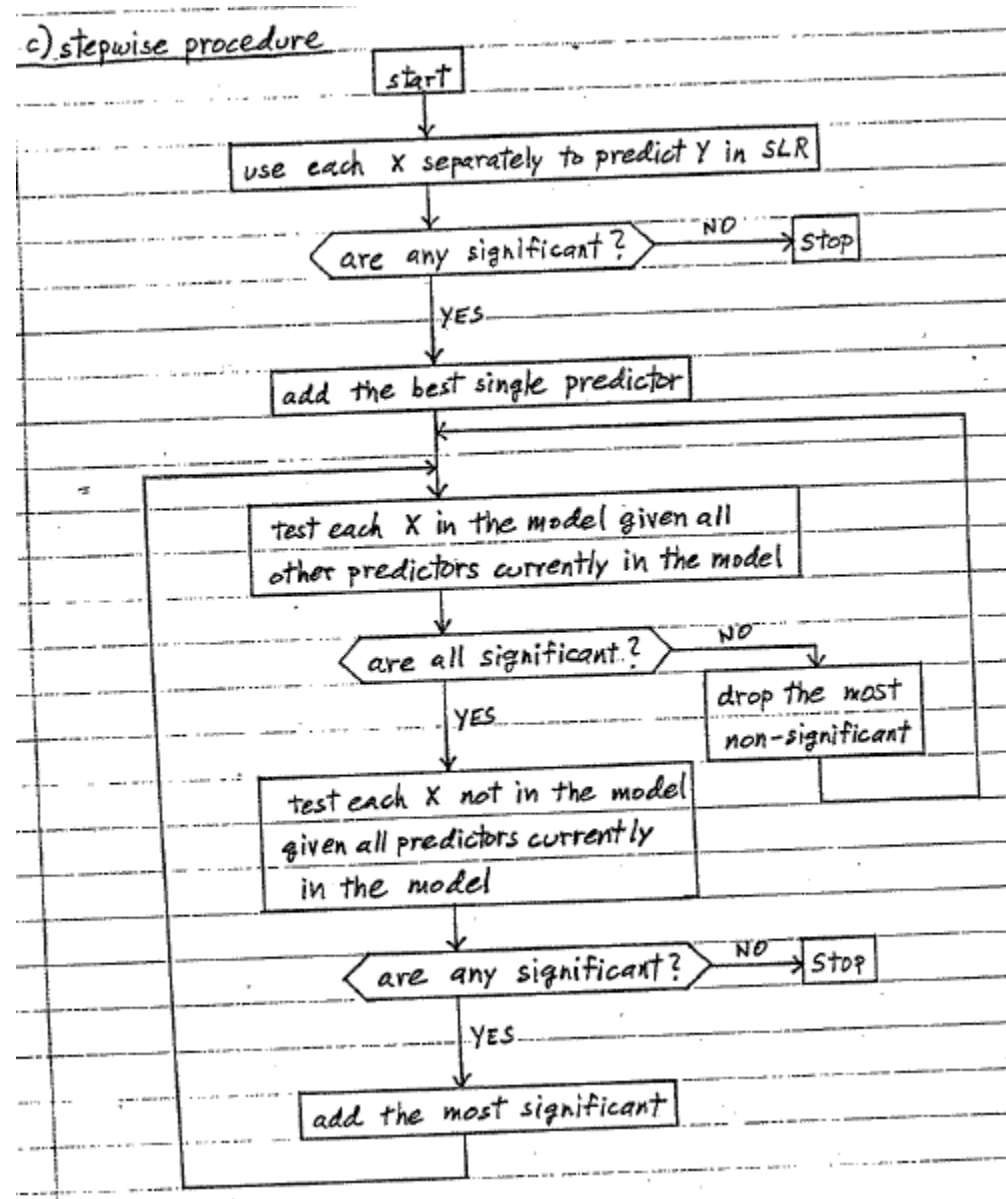
Draw the flow chart for the Backward Elimination:

Result



Draw the flow chart for the Stepwise Procedure:

Result



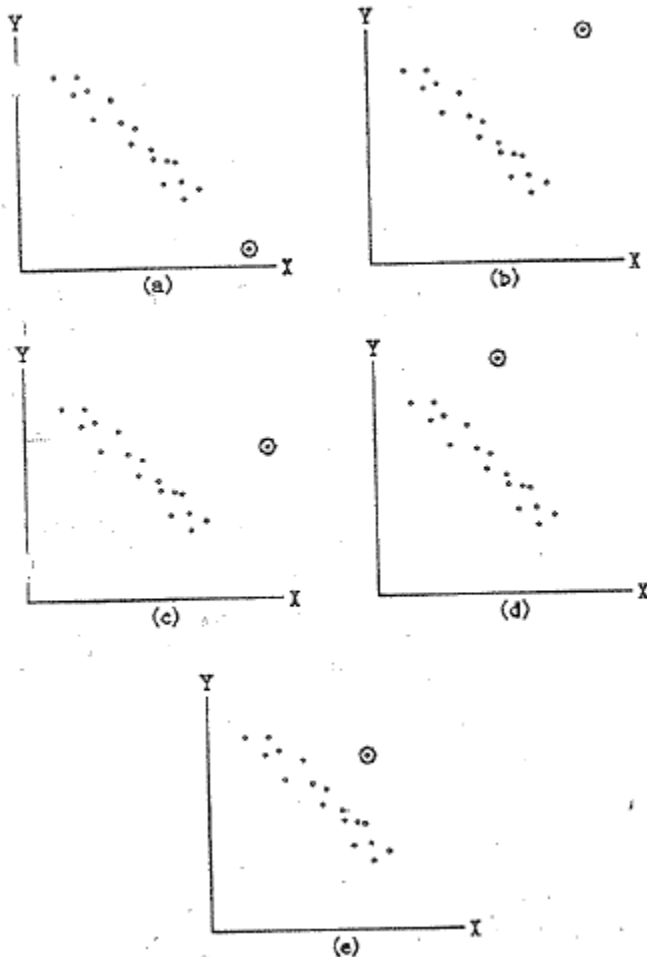
Explain the cross-validation procedure step by step.

Result

1. **Split Data:** Randomly divide data into training and validation sets.
2. **Fit Model:** Use the training set to select variables and fit the model.
3. **Predict:** Apply the model to the validation set to predict outcomes.

4. **Compare:** Check prediction accuracy by comparing predicted and actual values, often using sum of squared errors (SSE).
5. **Refit or Revise:** If the model performs well, refit using all data for stable estimates. If not, reduce predictors or revise the model to avoid overfitting.

Look at the Graps, and write the cause of outlier pattern or position (X or Y)



Result

- a. Location X and Y

- b. Location X and Y, and Pattern
- c. Location X and Pattern
- d. Location Y and Pattern
- e. Only Pattern