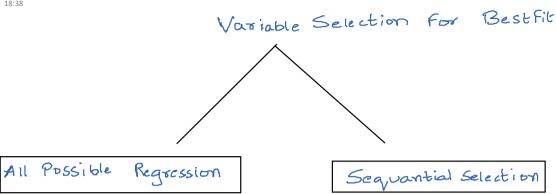
30 November 2017



All possible Ryression:

we need to consider all regression equations involving zero regressors in Y= Bo+E

-) Assume that, there are (K-1) Regressors (Regressors means x'rafiables)

-> As part of the model development we need to estimate 'k' parameters.

like Bo, B, Bz ... Bk-1

-: Total number of regressor model is 2

# Of Regressors	# of Model
0	F-1 CD
1 2	F-1 CD 12-1C, K-1C2
•	
K-I	K-10 K-1
	2 ^{K-1}

En: if we have 4'x' variables, we have these following all prossible Regressor models.

No		Regression Model	R ²	MS res	AdjR ²	Ср
	1	y=β0+ε	0	226.3	0	442.92
	2	y=β0+β1X1+ε	53.4	115.06	49.2	202.55
	3	y=β0+β2X2+ε	66.6	82.39	63.6	142.49
	4	y=β0+β3X3+ ε	28.6	176.3	22.1	315.16
	5	y=β0+β4Χ4+ε	67.5	80.35	64.5	138.73
	6	y=β0+β1X1+β2X2+ε	97.9	5.7	97.5	2.68
	7	y=β0+β1X1+β3X3+ε	54.8	122.7	45.8	198.1
	8	y=β0+β1X1+β4X4+ε	97.2	7.47	96.7	5.5
	9	y=β0+β2X2+β3X3+ε	84.7	41.54	81.7	62.44
	10	y=β0+β2X2+β4X4+ε	68	86.88	61.2	38.23
	11	y=β0+β3X3+β4X4+ε	93.5	17.57	92.2	22.37
	12	y=β0+β1X1+β2X2+β3X3+ε	98.2	5.34	97.6	3.04
	13	y=β0+β1X1+β2X2+β4X4+ε	98.2	5.33	97.6	3.02
	14	y=β0+β1X1+β3X3+β4X4+ε	98.1	5.64	97.5	3.5
	15	y=β0+β2X2+β3X3+β4X4+ε	97.3	8.2	96.3	7.34
	16	y=β0+β1X1+β2X2+β3X3+β4X4+ε	98.2	5.9	97.3	5

$$2^{(k-1)} = 2^{(k-1)}$$
= 16
= 16
 $2^{(k-1)} = 2^{(k-1)}$
= 24
 $2^{(k-1)} = 2^{(k-1)}$
= 24
 $2^{(k-1)} = 2^{(k-1)}$
= 26
 $2^{(k-1)} = 2^{(k-1)}$
= 27
 $2^{(k-1)} = 2^{(k-1)}$
= 26
 $2^{(k-1)} = 2^{(k-1)}$
= 27
 $2^{(k-1)} = 2^{(k-1)}$
= 26
 $2^{(k-1)} = 2^{(k-1)}$
= 26
 $2^{(k-1)} = 2^{(k-1)}$
= 27
 $2^{(k-1)} = 2^{(k-1)}$
= 28
 $2^{(k-1)} = 2^{(k-1)}$
= 26
 $2^{(k-1)} = 2^{(k-1)}$
= 26
 $2^{(k-1)} = 2^{(k-1)}$
= 27
 $2^{(k-1)} = 2^{(k-1)}$
= 28
 $2^{(k-1)} = 2^{(k-1)}$

-; These (2^{k-1}) equations are evaluated by below Conteria.

- 6 R2
- (2) Adjusted R2
- 3 MS Residual
- @ Marion Statistic (CP)

Mallow Statistic [Cp]:

Mallow statistic Measures the overall bias (4) Mean Square Error in the fitted model.

> = predicted Response

E(y) = True Response, 'K'is # of perameters (i.e. # of 3's)

 $C_{p} = \frac{\sum E(\hat{Y} - E(Y))}{\sigma^{2}}, \text{ Here we are dividing M.S.E with its S.D}$ $\Rightarrow \text{ Standadized Mean Square Error.}$

Here, if -E(Y) is nothing but (predicted value - population value) ie.

So, In Realtime we cannot find population value le E(y) so, we can estimated by

Ms (Full)

where Ke # of parameters including B.

Ms (full) = This is always for the full model

SSRes (16) = This is for k' parameters (i.e, (k-1) regressors)

 $Cn: Cp = \frac{SS_{Res}(x_1)}{M_{Ses}(x_1,x_2,x_3,x_4)} - 13 + 2(2)$

= 202.55

But the question is howdoyou find Best-Cp from all calculabed...! If any calculated Cp value is equilvalent to the Molyparameters i.e. going to be best Cp.