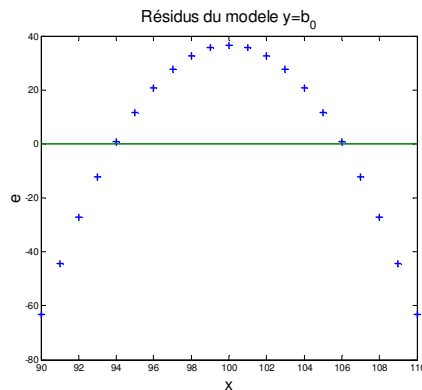


Sélection de variables

Cas 1: $Y=(X-100)^2$

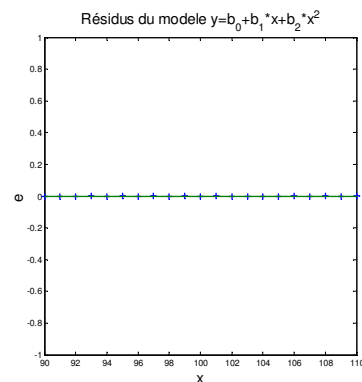
a) Sélection avant et « stepwise »:

La variable la plus corrélée est x^2 , mais $SCR_m=16.3$ et $SCE=22416$. L'ajout de la variable x^2 n'est donc pas significatif puisque $(16.3/1)/(22416/19)=.0138$, ce qui est très inférieur à une $F_{1,19,.05}=4.38$. Donc, aucune variable n'est sélectionnée mais les résidus sont inacceptables car ils montrent clairement qu'une tendance quadratique a été oubliée dans le modèle.



Supposons que l'on force le modèle $Y=b_0+b_1X+b_2X^2+e$.

On obtient les résidus suivants :



b) Élimination arrière : aucune variable ne peut être éliminée -> on a le bon modèle.

c) Statistique C_p :

- avec x et x^2 : $C_p=3$
- avec x , $C_p=1.1*10^{28}$
- avec x^2 , $C_p=1.1*10^{28}$

Clairement seul le modèle avec x et x^2 est adéquat, donc on a le bon modèle.

2^e cas : $Y=X^3+300*x+N(0,1)$

a) Sélection avant :

Étape	Variable	F ajout	F table (.05)
1	x^2	8080	4.38
2	x	6.5	4.41
3	x^3	58000	4.45

Les 3 ajouts sont significatifs. Le R^2 vaut 1. Toutefois, l'examen des coefficients et de leurs écarts-types révèle que le coefficient associé à x^2 et à la constante n'est pas significativement différent de zéro :

Variable	b	Écart-type
x	300.2	0.34
x ²	-.029	0.04
x ³	1.00	0.001
cte	-0.44	0.764

En reprenant la régression sans la constante et sans le x², on obtient à nouveau un R² de 1, et les coefficients de x et x³ sont fortement significatifs.

b) Élimination arrière :

SCE avec cte, x, x², x³: 18.23

SCE avec cte, x, x³: 18.80

SCE avec x, x³: 18.81

On peut éliminer x² car $\frac{(18.80 - 18.23)/1}{18.23/17} = 0.53 < F_{1,17,.05} = 4.45$.

On peut éliminer la constante, car : $\frac{(18.81 - 18.80)/1}{18.80/18} = 0.01 < F_{1,18,.05} = 4.41$

c) Stepwise : on inclut x², x, x³ puis on élimine x² et on élimine la constante.

d) Statistique C_p :

Modèle	C _p	Modèle	C _p
cte, x, x ² , x ³	4	x, x ² , x ³	3
cte, x, x ²	5.8*10 ⁵	x, x ²	9.2*10 ⁵
cte, x, x ³	2.5	x, x ³	1.2
cte, x ² , x ³	7.8*10 ⁵	x ² , x ³	8*10 ⁵
cte, x ¹	1.9*10 ⁷	x ¹	3.4*10 ⁷
cte, x ²	7.9*10 ⁵	x ²	3.9*10 ⁶
cte, x ³	1.1*10 ⁷	x ³	3.9*10 ⁷

Les seuls modèles acceptables sont (cte, x, x² et x³), (cte, x, x³), (x, x², x³) et (x, x³), celui-ci étant le meilleur de tous au sens de la statistique C_p. Si l'on examine le modèle cte, x, x², x³, on constate que la constante et le coefficient pour b² ne sont pas significatifs puisque l'intervalle de confiance inclut le 0 au niveau α=0.05. (b₀/s_{b0}=-0.57; b₁/s_{b1}=885; b₂/s_{b2}=-0.7; b₃/s_{b3}=763).

Note : Les coefficients trouvés avec x et x³ (modèle sans constante) sont b₁=299.9 et b₃=1.00, soit presque exactement le modèle théorique. Le R² vaut presque 1.

Les résidus ne montrent pas d'anomalies particulières :

