

UNIVERSITÉ CLAUDE BERNARD LYON I

MASTER DATA SCIENCE

MODÈLES DE RÉGRESSION

Régression des données d'extraction de puits de pétrole au Canada

Auteurs :

Bruno DUMAS

Jules SAUVINET

Professeur :

François WAHL

14 janvier 2017



Résumé

L'objectif est d'effectuer des régressions des valeurs d'extractions des puits de pétrole du Canada afin de pouvoir prédire les futures données d'extraction. Les courbes des valeurs des données d'extraction donnent une classification de la qualité des puits. L'étiquetage de la qualité de chaque puits à été effectué par des experts. L'objectif de ce travail est de mettre en place une classification automatique de ces puits : La démarche proposée est décrite ci-dessous. L'idée est de remplacer ces courbes par des fonctions paramétriques. Pour cela, plusieurs régressions vont être envisagées afin d'avoir la meilleure prédiction/classification possible.

Table des matières

1	Régressions polynomiales	2
2	Régressions exponentielles	3
3	Courbes hautes et basses à 95%	5
4	Reclassement avec régression logistique	6
5	Gestion des spikes et lissage des courbes	8
A	Code de la question 1	11
B	Code de la question 2	12
C	Code de la question 3	14
D	Code de la question 4	16
E	Code de la question 5	18

1 Régressions polynomiales

Une façon simple est d'ajuster un polynôme de degré faible sur chacune des courbes et de voir si les coefficients présentent des clusters, c'est à dire des groupes de points distincts quand on les regarde dans l'espace. On essaiera des polynômes de degré 0, 1, 2, 3, et 4. On présentera les courbes de production simulées obtenues, comme dans la figure ci-dessous.

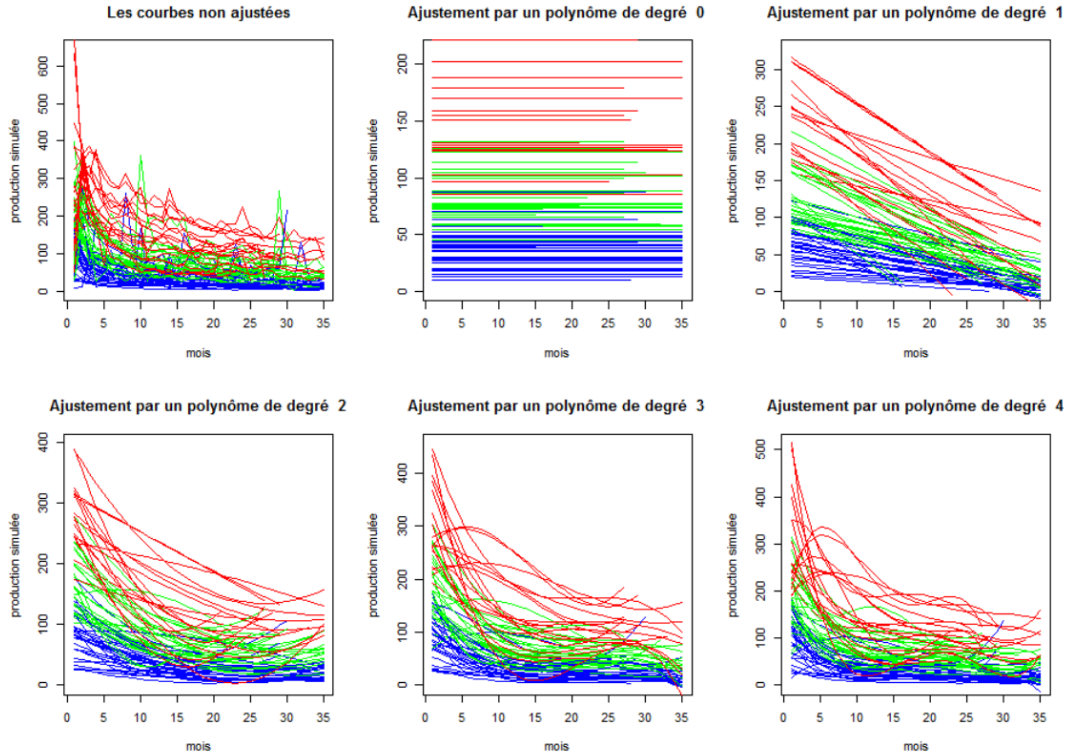


FIGURE 1 – Courbes non régressées, et courbes de régressions polynomiales de degrés 0,1,2,3 et 4

La moyenne des R-Squared ajustés pour les 75 courbes pour les 5 types de régression sont respectivement :

0.0000000, 0.4853609, 0.6622269, 0.7202006 et 0.7557633.

On voit que plus le degré du polynôme augmente, plus la régression est de qualité au critère du R-Squared (les valeurs prédites se rapprochent des valeurs des observations).

Si on se réfère à la figure ci-dessus, on observe que la plupart des simulations présente une remontée au bout de quelques mois, que certaines simulations sont concaves au lieu d'être convexes, voire que certaines d'entre elles pourraient avoir des valeurs négatives (autour des 35 mois d'exploitation).

2 Régressions exponentielles

Une idée simple pour corriger ces défauts est d'utiliser une autre forme paramétrique pour les simulations. Une suggestion immédiate pour qui a un peu l'habitude de ces courbes est une forme exponentielle du style : $y = k_0.e^{-k_1.t}$ où y est la production, t le mois, et k_0 et k_1 deux paramètres à déterminer.

Régression exponentielle $\log(y) = ax + b$

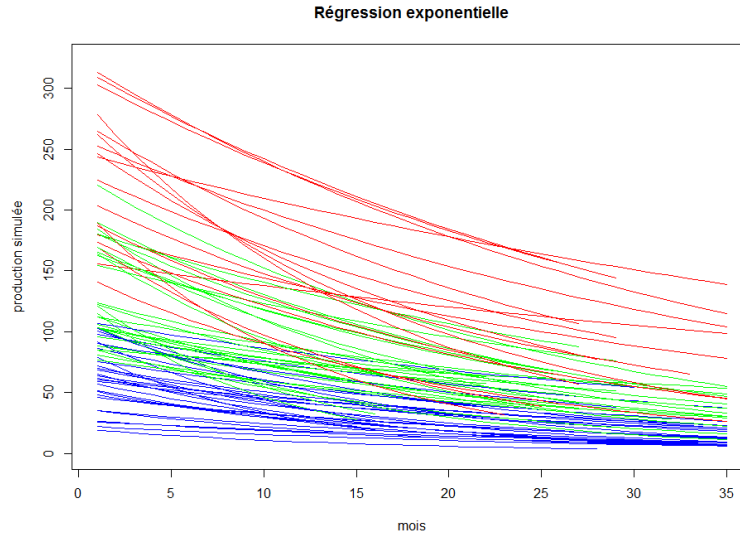


FIGURE 2 – Courbes ajustées avec une fonction exponentielle et l'outil lm de R

Régression exponentielle $y = k_0.e^{-k_1.x}$

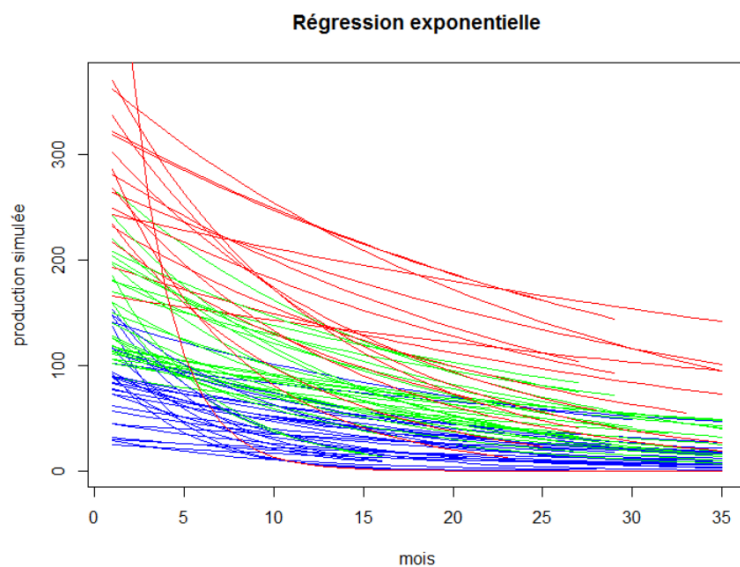


FIGURE 3 – Courbes ajustées avec une fonction exponentielle et l'outil nls de R

La régression $\log(y) = ax+b$ faite avec `lm` est plus fidèle aux données, on s'intéressera donc à celle-ci.

La moyenne du "adjusted R-Squared" pour les 75 régressions exponentielles est cette fois-ci de 0.6754964, donc de moins bonne qualité à priori que les régressions polynomiales de degré supérieures ou égale à 3.

Néanmoins, quand on s'intéresse au détail des R-Squared, on trouve de très bonne régressions avec un R-Squared très bon comme des régressions avec un R-Squared très mauvais. Cela est probablement dû à certaines valeurs "outliers" qui ne permettent pas à la fonction exponentielle de s'adapter aux courbes.

La partie 5 et le lissage des courbes permettra de pallier ce problème et probablement de montrer que la régression exponentielle est adaptée si un travail d'atténuation des "pics" est fait au préalable.

10 premiers R-Squared :

0.8011383, 0.4182548, 0.7403922, 0.7585257, 0.4186529, 0.6888943, 0.3297291, 0.818889, 0.7745934 et 0.7527809.

On peut d'ores et déjà constater que les trois classes de puits définies par les experts sont visuellement distinguables, bien qu'un peu entremêlées. Ceci indique qu'il sera possible de mettre en place un modèle de classification des puits.

Autres possibilités

Nous avons pensé à un ajustement par fonction inverse du type $y = \frac{a}{x} + b$ avec y la production, x le mois et a et b des constantes à déterminer. On peut aussi penser à une gaussienne, dont l'allure s'approche de celles des courbes.

3 Courbes hautes et basses à 95%

Quelles sont les incertitudes sur les régressions des points 2 ? Plus concrètement, on vous demande de tracer pour un exemple de chaque type de courbe, la courbe haute (à 95%) et la courbe basse (toujours à 95%)

Intervalles de confiance avec nls

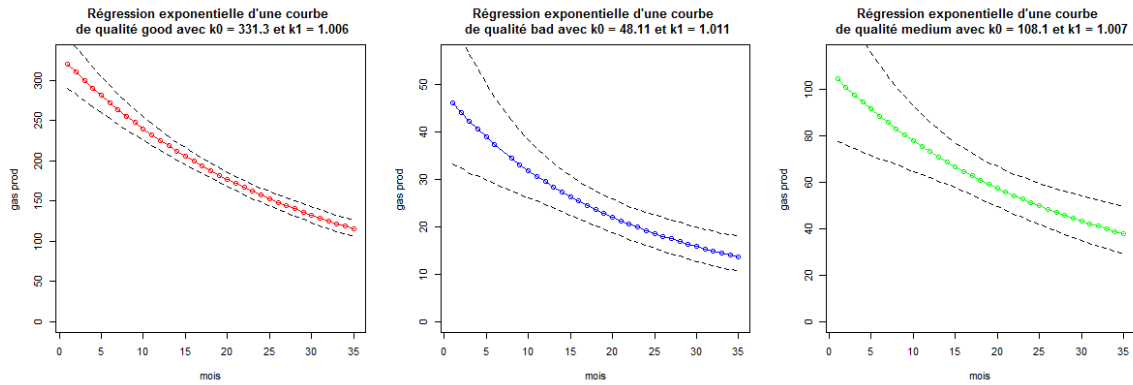


FIGURE 4 – Courbes hautes et basses à 95% pour un exemple de chaque classe de courbes

Intervalles de confiance avec lm

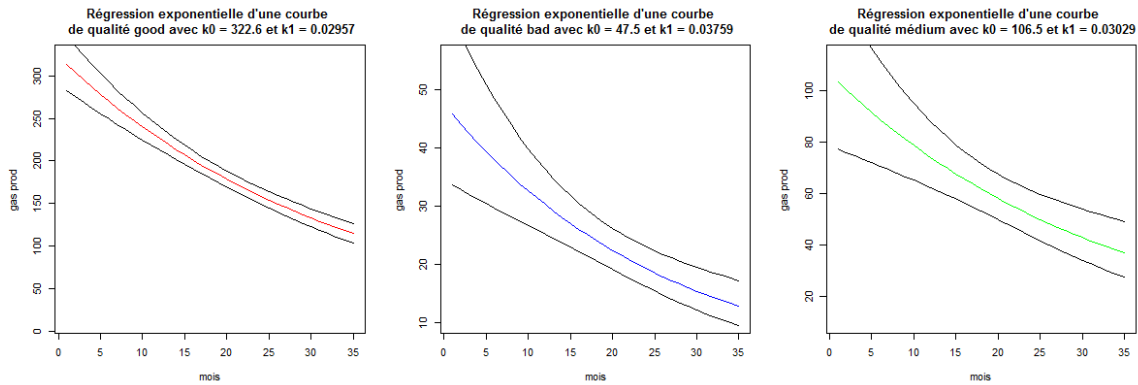


FIGURE 5 – Courbes hautes et basses à 95% pour un exemple de chaque classe de courbes

Comme le montre les figures ci-dessus, il y a davantage d'incertitude sur les valeurs des courbes de qualité 'bad' et 'médium'. Cela induit que les valeurs de ces courbes suivent des tendances plus compliquées à ajuster. En effet, les valeurs d'extraction pour ces puits sont probablement plus imprévisibles et suivent parfois quelques oscillations. En outre, les incertitudes aux valeurs d'extraction pour les premiers et derniers mois sont également plus fortes.

4 Reclassement avec régression logistique

En examinant le graphe $k1$ fonction de $k0$, on se rend compte que certaines courbes classées 'Good' par les experts donnent l'impression d'être plutôt 'medium', tandis que certaines 'bad' pourraient être aussi 'medium'.

Avez-vous des suggestions sur 5 courbes au plus qui pourraient être mal classées ?

Justifiez vos choix (i.e. une façon de faire est d'effectuer une régression logistique dont le y est la classe prédite par l'expert et les x sont les coefficients $k0$ et $k1$, et d'examiner comment la régression est améliorée en changeant la classe d'un point).

Clustering des courbes en fonction de $k0$ et $k1$ avant reclassement

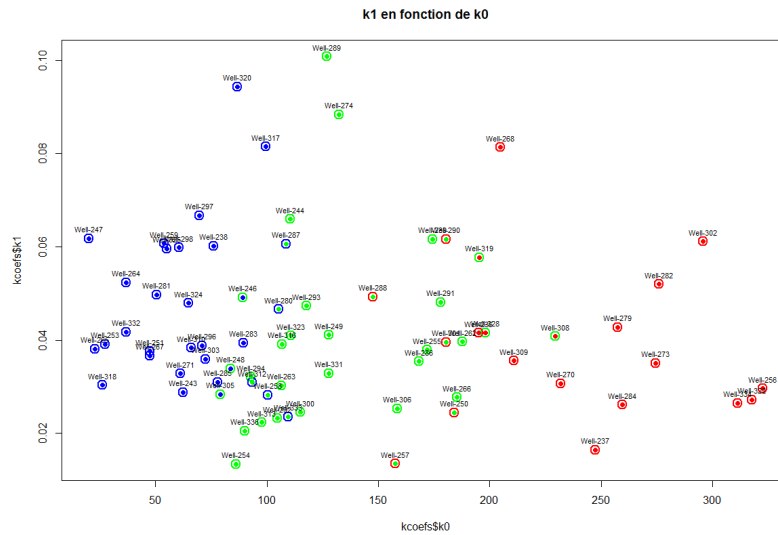


FIGURE 6 – Coefficients $k1$ en fonction de $k0$ et classes des courbes

La couleur du cercle extérieur est la classe donnée par les experts et le cercle intérieur est sa classe prédite par le classifieur.

16 courbes sont mal classées d'après les prédictions et l'AIC de la régression est de 74.63465.

On essaye d'améliorer le taux du classifieur en modifiant la classe prédite par les experts sur certaines courbes.

On change ainsi la classe de 5 puits.

On se limite au reclassement de 5 puits car on ne veut pas trop modifier le modèle de départ et laisser de la souplesse au classifieur si celui-ci doit traiter la classification de nouveaux puits à l'avenir.

On procède de manière itérative. On choisit un puits qui semble très éloigné du reste de sa classe et on change sa classe donnée par les experts pour celle prédite par le classifieur. On relance le modèle du classifieur sur cette nouvelle base de données. On compte les points mal classés avec ce nouveau classifieur. Si le puits choisi ne diminue pas le nombre de points bien classés, on lui réattribue sa valeur originelle et on choisit un autre point. Quand on trouve un point qui convient, on répète ces

opérations jusqu'à atteindre 5 points.

On aurait pu automatiser cette opération qui est un problème d'optimisation simple. Ce problème peut se résumer par la phrase suivante :

Trouver la séquence de 5 puits qui minimise le nombre de puits mal classés

Après interprétations graphiques et des tests de prédiction de classe avec ou sans changement de classe des 16 puits mal classés, on détermine les 5 puits qui réajustent au mieux le modèle et améliorent le classifieur :

Les courbes des puits 'Well - 288' de good à médium, 'Well - 333' de bad à médium, 'Well - 246' de médium à bad, 'Well - 257' de good à médium, et 'Well - 258' de bad à médium.

Clustering des courbes en fonction de $k0$ et $k1$ après reclassement

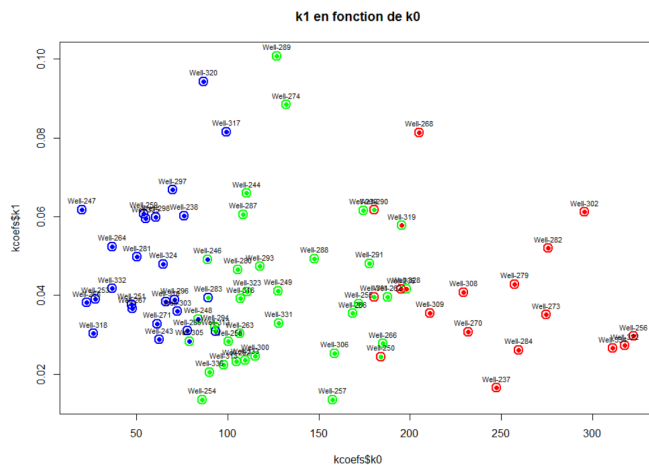


FIGURE 7 – Coefficients $k1$ en fonction de $k0$ et classes des courbes

L'AIC de la régression est de 52.52519 et il n'y a plus que 10 courbes mal classées pour 5 de changées sur les 16 de départ. Soit 1 courbe qui dont la prédiction de classe par le classifieur s'est accordée sur sa vraie classe.

A présent que nous avons établi un classifieur, nous pouvons nous demander s'il est possible d'en améliorer la précision. On peut par exemple tenter de produire un meilleur modèle sur les données en lissant les courbes avant d'effectuer la régression.

5 Gestion des spikes et lissage des courbes

Régression polynomiale de degré 3 avec smooth

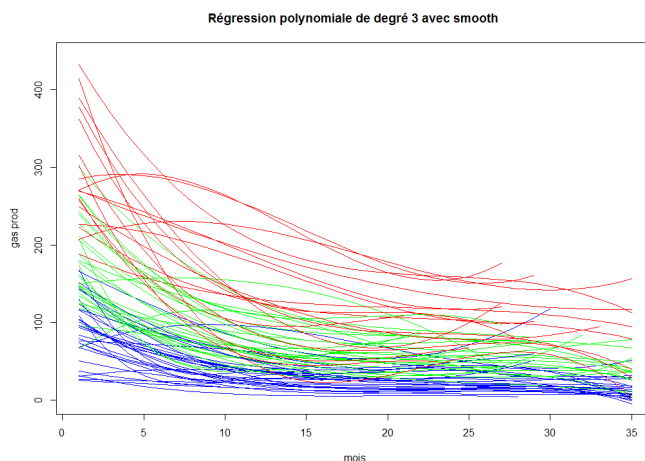


FIGURE 8 – Régression polynomiale de degré 3 avec courbes lissées au préalable avec loess

On obtient une moyenne de R-Squared de 0.9802724 cette fois-ci après lissage des courbes. La régression polynomiale devient ainsi très performante. Il reste à savoir si l'on peut se permettre l'approximation faite par le lissage des courbes et si les valeurs induites par les spikes avaient une pertinence intransigible.

Toutefois, la régression avec smooth ne règle pas les valeurs qui remontent, les simulations concaves au lieu d'être convexes, et les potentielles valeurs négatives autour des 35 mois.

Régression exponentielle avec lissage des spikes

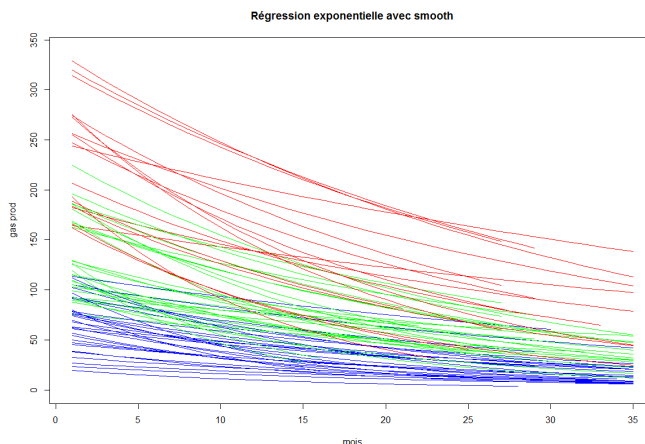


FIGURE 9 – Régression exponentielle lm avec courbes lissée au préalable avec loess

On obtient une moyenne de R-Squared de 0.7569813 cette fois-ci après lissage des courbes. La régression est sensiblement améliorée avec le lissage des courbes mais reste moins efficace que la régression polynomiale de degré 3.

Reclassement après lissage des courbes et ajustement exponentiel

On effectue à nouveau une régression logistique à partir des courbes lissées comme à la question 4. On obtient cette fois-ci toujours 16 courbes mal classées, un AIC de 85 et le graphique suivant :

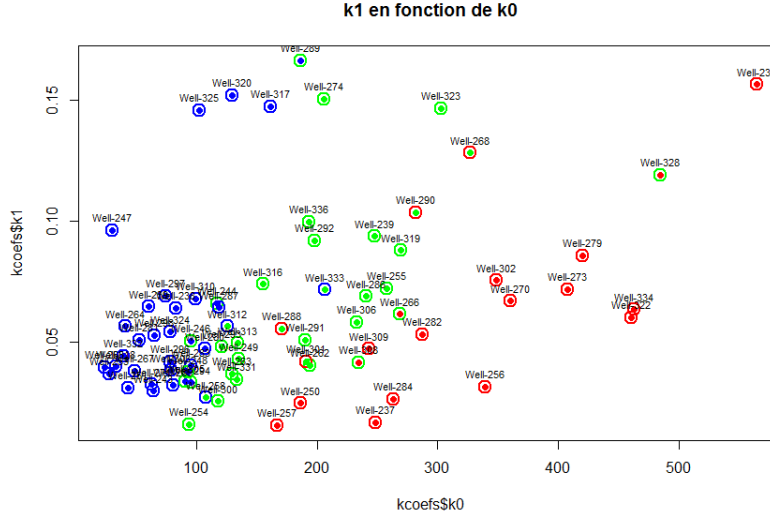


FIGURE 10 – Clustering sur modèle exponentiel avec courbes lissées

On se limite comme dans la partie 3 à seulement 5 reclassement de courbes. Après reclassement des puits 'Well - 290', 'Well - 333', 'Well - 312', 'Well - 288', 'Well - 258', on obtient un AIC de 56 et plus que 8 courbes mal classées en faisant une nouvelle régression logistique, soit 3 courbes dont la classe prédite à cette fois-ci été égale à celle observée.

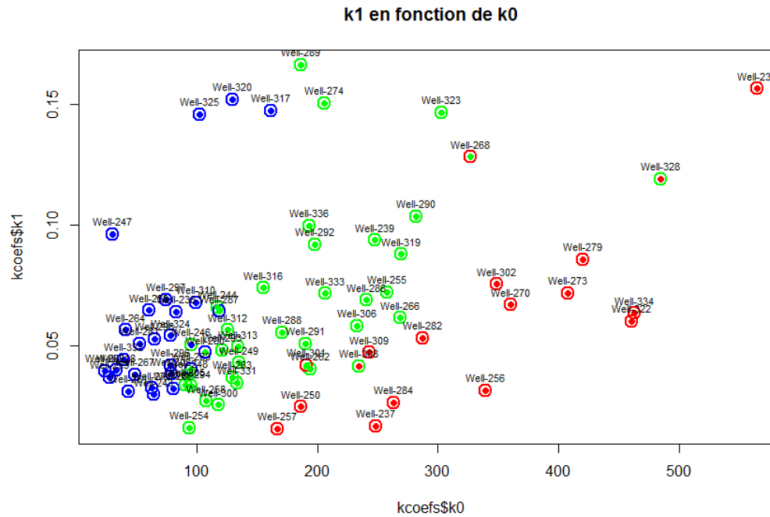


FIGURE 11 – Clustering sur modèle exponentiel avec courbes lissée après reclassement

Reclassement après lissage des courbes et ajustement polynomial

On effectue une régression logistique à partir des courbes lissées et ajuster par un polynôme de degré 3 cette fois-ci. On obtient cette fois-ci toujours 12 courbes mal classées, un AIC de 78.25367 et le graphique de dépendance des 4 coefficients entre eux suivant :

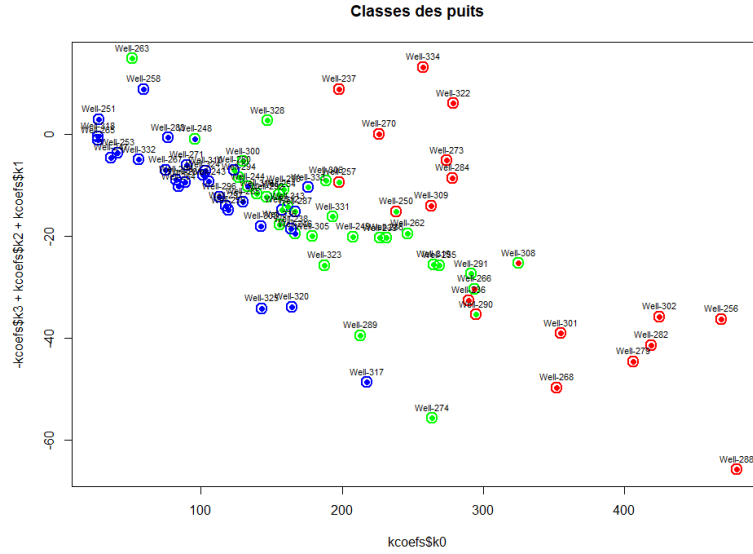


FIGURE 12 – Clustering sur modèle polynomial de degré 3 avec courbes lissées

Après reclassement des puits 'Well - 257', 'Well - 250', 'Well - 333', 'Well - 266' on obtient un AIC de 64.16384 et plus que 6 courbes mal classées en faisant une nouvelle régression logistique, soit 2 courbes dont la classe prédite à cette fois-ci été égale à celle observée.

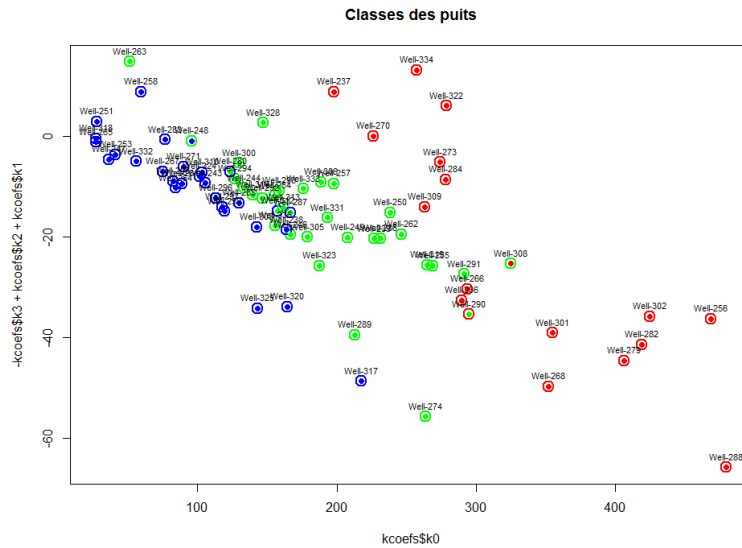


FIGURE 13 – Clustering sur modèle polynomial de degré 3 avec courbes lissées et reclassement

Acknowledgments

Nos remerciements vont à François Wahl pour l'enseignement de son cours "Modèles de régression" à l'Université Claude Bernard Lyon I et son accompagnement durant ce projet.

A Code de la question 1

```
15 #récupération des données du fichier EXCEL
16 perl <- 'C:\\Strawberry\\perl\\bin\\perl.exe'
17 datapuits = read.xls(file.path("data/FW_Donnees_Puits.xlsx"), perl=perl)
18
19 #on transpose les données pour avoir les puits en colonne et les mois en ligne
20 puits = setNames(data.frame(t(datapuits[, -1])), datapuits[, 1])
21 puits
22
23 #####
24 ## QUESTION 1 - Régressions polynomiales##
25 #####
26
27 par(mfrow=c(2,3))
28 seq1 <- 1:4 #on va tracer 6 graphiques
29 seq2 <- 1:75 #les 75 puits
30 r2 <- numeric(5) #la moyenne du rsquared
31 for (j in seq1){
32   r2j = 0.0 # le rsquared d'une regression
33   for (i in seq2){
34     mois <- 1:35 #l'abscisse
35     v <- as.vector(puits[,i]) #on récupère les données du puits i
36
37     col="black" #la couleur en fonction de la classification de qualité
38     if (v[36] == "Good"){col = "red"}else if (v[36] == "medium"){col = "green"}else {col = "blue"}
39
40     v <- as.numeric(v[1:35])
41
42     #on plot pas les 0 qui sont des ND
43     nd <- which(v %in% 0)
44     v <- v[v != 0]
45     mois <- mois[!mois %in% nd]
46
47     #tracé de la figure 1 : les données de production
48     if (j==1){
49       if (i==1){
50         plot(mois,v,type="l",col=col, ylab="production simulée", main="Les courbes non ajustées",ylim=c(0,max(v)+10))
51       }
52       else
53         lines(mois,v,type="l",col=col)
54     }
55     #tracé de la figure 2 : les courbes de production obtenues avec des polynômes de degré 0, 1, 2, 3, et 4
56     else{
57       if (j==0){
58         fit2 <- lm(v ~ 1)
59       }else {
60         fit2 <- lm(v ~ poly(mois, j, raw=TRUE))
61       }
62       if (i==1){
63         plot(mois, predict(fit2), type="l",col=col, ylab="production simulée", lwd=1, main=paste("Ajustement par un polynôme de c
64       )
65       lines(mois, predict(fit2), col=col, lwd=1)
66     }
67     if (j >= 0) {r2j=r2j+summary(fit2)$adj.r.squared}
68   }
69   #on calcule la moyenne des R-Squared pour chaque type de régression
70   if (j >= 0){
71     r2[j+1]=r2j/i
72   }
73 }
74
```

FIGURE 14 – Code de la question 1 des régressions polynomiales

B Code de la question 2

```
70 #####
71 ## QUESTION 2 - Régression exponentielle ##
72 #####
73
74 #Avec lm
75 rse2=0
76 r2e1 <- numeric(1)
77 for (i in seq2){
78   mois <- 1:35
79   v <- as.vector(puits[,i])
80
81   col="black"
82   if (v[36] == "Good"){
83     col = "red"
84   }else if (v[36] == "medium"){
85     col = "green"
86   }else
87     col = "blue"
88
89   v <- as.numeric(v[1:35])
90
91   nd <- which(v %in% 0)
92   v <- v[v != 0]
93   mois <- mois[!mois %in% nd]
94
95   expfit <- lm(log(v) ~ mois)
96
97   if (i==3){
98     #residus=v-exp(expfit$fitted.values)
99     #plot(v,residus, main=paste("Graphique des résidus de la régression exponentielle"), xlab="production simulée",
100
101     plot(mois,exp(predict(expfit)),type="l",col=col, ylab="production simulée",
102          main=paste("Régression exponentielle"),ylim=c(0,max(exp(predict(expfit)))+10))
103
104   }
105   else {
106     lines(mois,exp(predict(expfit)),type="l",col=col)
107   }
108   rse2 = rse2 + exp(sigma(expfit))
109 }
110 rse2 = rse2 / 75|
111 rse2
112 summary(expfit)
```

FIGURE 15 – Code de la question 2 des régressions exponentielles, partie 1

```

124 #Avec nls
125 rse = 0
126 seq2 <- 1:75
127 r2e2 <- numeric(1)
128 for (i in seq2){
129   mois <- 1:35
130   v <- as.vector(puits[,i])
131
132   col="black"
133   if (v[36] == "Good"){
134     col = "red"
135   }else if (v[36] == "medium"){
136     col = "green"
137   }else {
138     col = "blue"
139   }
140
141   v <- as.numeric(v[1:35])
142
143   #on plot pas les 0 qui sont des ND (d'apres moi)
144   nd <- which(v %in% 0)
145   v <- v[v != 0]
146   mois <- mois[!mois %in% nd]
147   df <- data.frame(mois, v)
148   df$lv <- df$v
149
150   k0start=-300
151   k1start=0.1
152
153   m <- nls(lv ~ k0*exp(-k1*mois), start=c(k0=k0start, k1=k1start), df)
154   summary(m)
155
156   if (i==1){
157     plot(df$mois,predict(m),type="l",col=col, ylab="production simulée", xlab="mois", main=paste("Régression expone
158   })
159   lines(df$mois,predict(m),type="l",col=col)
160
161   rse = rse + sigma(m)
162 }
163 rse = rse / 75
164 rse
165 summary(m)
166

```

FIGURE 16 – Code de la question 2 des régressions exponentielles, partie 2

C Code de la question 3

```

169 #####
170 ## QUESTION 3 Courbe haute + courbe basse à 95% ##
171 #####
172
173 #avec nls + predict_NLS
174 plotGood <- TRUE
175 plotMed <- TRUE
176 plotBad <- TRUE
177 for (i in seq2){
178   mois <- 1:35
179   v <- as.vector(puits[,i])
180   classif = v[36]
181
182   v <- as.numeric(v[1:35])
183
184   nd <- which(v %in% 0)
185   v <- v[v != 0]
186   mois <- mois[!mois %in% nd]
187   df <- data.frame(mois, v)
188   df$lv <- log(df$v)
189
190   col="black"
191
192   k0start=400
193   k1start=0.01
194
195   m <- nls(lv ~ k0*exp(-k1*mois), start=c(k0=k0start, k1=k1start), df)
196   summary(m)
197
198   k0 = signif(exp(coef(m)[1][["k0"]]), digits = 4)
199   k1 = signif(exp(coef(m)[2][["k1"]]), digits = 4)
200
201
202   if (classif == "Good" && plotGood == TRUE){
203     predict1 = predictNLS(m, df)
204     plotGood = FALSE
205     col = "red"
206     plot(df$mois,exp(predict(m)),type="p",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle"))
207     lines(df$mois,exp(predict(m)),type="l",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle"))
208
209     lines(mois,exp(predict1[,6]),type="l",col="black",lwd = 1,lty=2)
210     lines(mois,exp(predict1[,7]),type="l",col="black",lwd = 1,lty=2) |
211
212   }else if (classif == "medium" && plotMed == TRUE){
213     predict2 = predictNLS(m, df)
214     plotMed = FALSE
215     col = "green"
216     plot(df$mois,exp(predict(m)),type="p",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle"))
217     lines(df$mois,exp(predict(m)),type="l",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle"))
218
219     lines(mois,exp(predict2[,6]),type="l",col="black",lwd = 1,lty=2)
220     lines(mois,exp(predict2[,7]),type="l",col="black",lwd = 1,lty=2)
221
222   }else if (classif == "bad" && plotBad == TRUE) {
223     predict3 = predictNLS(m, df)
224     plotBad = FALSE
225     col = "blue"
226     plot(df$mois,exp(predict(m)),type="p",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle"))
227     lines(df$mois,exp(predict(m)),type="l",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle"))
228
229     lines(mois,exp(predict3[,6]),type="l",col="black",lwd = 1,lty=2)
230     lines(mois,exp(predict3[,7]),type="l",col="black",lwd = 1,lty=2)
231   }
232 }
233

```

FIGURE 17 – Code de la question 3 des intervalles de confiance, partie 1

```

234 #avec predict
235 seq1 <- 0:4
236 #les 75 premiers puits
237 seq2 <- 1:75
238
239 plotGood <- TRUE
240 plotMed <- TRUE
241 plotBad <- TRUE
242
243 par(mfrow=c(2,3))
244 for (i in seq2){
245   mois <- 1:35
246   v <- as.vector(puits[,i])
247   classif = v[36]
248
249   col="black"
250
251   v <- as.numeric(v[1:35])
252
253   nd <- which(v %in% 0)
254   v <- v[v != 0]
255   mois <- mois[!mois %in% nd]
256
257   expfit <- lm(log(v) ~ mois)
258
259   k0start = signif(exp(expfit$coefficients[1]), digit=4)
260   k1start = signif(-expfit$coefficients[2], digit=4)
261
262   newdat <- data.frame(mois)
263
264   if (classif == "Good" && plotGood == TRUE){
265     plotGood = FALSE
266     col = "red"
267
268     predG = predict(expfit, newdat, interval="confidence", level=0.95)
269     epredG = exp(predG)
270
271     plot(mois,epredG[,1],type="l",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle d'une"))
272
273     lines(mois,epredG[,2],type="l",col="black",ylim=c(0,max(exp(predict(expfit)))+10))
274     lines(mois,epredG[,3],type="l",col="black",ylim=c(0,max(exp(predict(expfit)))+10))
275   }else if (classif == "medium" && plotMed == TRUE){
276     plotMed = FALSE
277     col = "green"
278
279     predG = predict(expfit, newdat, interval="confidence", level=0.95)
280     epredG = exp(predG)
281
282     plot(mois,epredG[,1],type="l",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle d'une"))
283
284     lines(mois,epredG[,2],type="l",col="black",ylim=c(0,max(exp(predict(expfit)))+10))
285     lines(mois,epredG[,3],type="l",col="black",ylim=c(0,max(exp(predict(expfit)))+10))
286   }else if (classif == "bad" && plotBad == TRUE) {
287     plotBad = FALSE
288     col = "blue"
289
290     predG = predict(expfit, newdat, interval="confidence", level=0.95)
291     epredG = exp(predG)
292
293     plot(mois,epredG[,1],type="l",col=col, ylab="gas prod", xlab="mois", main=paste("Régression exponentielle d'une"))
294
295     lines(mois,epredG[,2],type="l",col="black",ylim=c(0,max(exp(predict(expfit)))+10))
296     lines(mois,epredG[,3],type="l",col="black",ylim=c(0,max(exp(predict(expfit)))+10))
297   }
298 }

```

FIGURE 18 – Code de la question 3 des intervalles de confiance, partie 2

D Code de la question 4

```
301 #####
302 ## QUESTION 4 Suggestions sur 5 courbes mal classées ##
303 #####
304
305 improvePredict <- function(colorsPred = c(), badClass = c()){
306   #On stocke les valeurs de k0, k1 et la classe (couleur)
307   k0s <- numeric(75)
308   k1s <- numeric(75)
309   colors <- numeric(75)
310   par(mfrow=c(1,1))
311   r2el <- numeric(1)
312   for(i in seq2){
313     mois <- 1:35
314     v <- as.vector(puits[,i])
315
316     col="black"
317     if (v[36] == "Good"){
318       if(i %in% badClass){
319         col = kcoefs$colPred[i]
320       } else {
321         col = "red"
322       }
323     } else if (v[36] == "medium"){
324       if(i %in% badClass){
325         col = kcoefs$colPred[i]
326       } else {
327         col = "green"
328       }
329     } else {
330       if(i %in% badClass){
331         col = kcoefs$colPred[i]
332       } else {
333         col = "blue"
334       }
335     }
336
337     v <- as.numeric(v[1:35])
338
339     nd <- which(v %in% 0)
340     v <- v[v != 0]
341     mois <- mois[!mois %in% nd]
342
343     expfit <- lm(log(v) ~ mois)
344     k0i = exp(expfit$coefficients[1])
345     k1i = -expfit$coefficients[2]
346     k0s[i] = k0i
347     k1s[i] = k1i
348     colors[i] = col
349
350
351     rsquared = summary(expfit)$adj.r.squared
352     r2el=r2el+summary(expfit)$adj.r.squared
353   }
354
355   r2el=r2el/i
356   kcoefs <- c()
357   kcoefs$k0 <- k0s
358   kcoefs$k1 <- k1s
359   kcoefs$col <- colors
360
361   clustering <- multinom(col ~ k0 + k1, data = kcoefs)
362   summary(clustering)
363
364   names <- lapply(datapuits$V1, as.character)
365
366   kcoefs$names <- names
367   kcoefs$model <- clustering
368
369   kcoefs$colPred <- clustering$lev[predict(clustering)]
370
371   # Compter les courbes mal classées
372   count=0
373   if(length(kcoefs$colPred) > 0){
374     for(i in 1:75){
375       if (kcoefs$col[i] != kcoefs$colPred[i]){count=count + 1}
376     }
377   }
```

FIGURE 19 – Code de la question 4 des reclassements de courbes après régression logistique, partie 1

```

378 kcoefs$badPredictCount <- count
379 kcoefs$correctedPoints <- length(badClass)
380
381 plot(kcoefs$k0,kcoefs$k1,type="p",pch=1,cex=2, lwd = 2,col=kcoefs$col, main="k1 en fonction de k0")
382 lines(kcoefs$k0,kcoefs$k1,type="p",pch=19,cex=1,col=clustering$lev[predict(clustering)],main="k1 en fonction de k0")
383 text(kcoefs$k0, kcoefs$k1, labels=names, cex= 0.7, pos=3)
384
385 print(kcoefs$badPredictCount)
386 print(kcoefs$correctedPoints)
387
388 return(kcoefs)
389 }
390
391 removePoint <- function(badClass, pointNumber){
392   badClass=c(badClass, which(kcoefs$names == paste("well", pointNumber, sep="-")))
393   kcoefs = improvePredict(colorsPred = kcoefs$colPred,badClass = badClass)
394   return (badClass)
395 }
396
397 badClass=c()
398 kcoefs = improvePredict()
399
400 # On a 16 courbes mal prédites. On en sélectionne 5.
401
402 # Courbe n°75
403 badClass = removePoint(badClass, 288)
404 # Courbe n°47
405 badClass = removePoint(badClass, 333)
406 # Courbe n°39
407 badClass = removePoint(badClass, 246)
408 # Courbe n°53
409 badClass = removePoint(badClass, 257)
410 # Courbe n°11
411 badClass = removePoint(badClass, 258)
412

```

FIGURE 20 – Code de la question 4 des reclassements de courbes après régression logistique, partie 2

E Code de la question 5

```
416 #####
417 ## QUESTION 5 - Gestion des spikes (smoothing curves) ##
418 #####
419
420 #fait avec loess
421 seq1 <- 0:4
422 #les 75 premiers puits
423 seq2 <- 1:75
424 #polynomial de degré 3
425 r2p3 <- numeric(1)
426 for (i in seq2){
427   mois <- 1:35
428   v <- as.vector(puits[,i])
429
430   #la couleur en fonction de la classification de qualité
431   col="black"
432   if (v[36] == "Good"){
433     col = "red"
434   }else if (v[36] == "medium"){
435     col = "green"
436   }else
437     col = "blue"
438
439   v <- as.numeric(v[1:35])
440
441   nd <- which(v %in% 0)
442   v <- v[v != 0]
443   mois <- mois[!mois %in% nd]
444
445   smooth <- loess(v~mois)
446   fit3 <- lm(smooth$fitted ~ poly(mois, 3, raw=TRUE))
447
448   if (i==1){
449     plot(mois, predict(fit3), type="l", col=col, ylab="gas prod", lwd=1, main="Régression polynomiale de degré 3 avec")
450   } else {
451     lines(mois, predict(fit3), col=col, ylab="gas prod", lwd=1)
452   }
453   lines(mois, predict(fit3), col=col, ylab="gas prod", lwd=1)
454
455   rsquared = summary(fit3)$adj.r.squared
456   r2p3=r2p3+summary(fit3)$adj.r.squared
457 }
458
```

FIGURE 21 – Code de la question 5 des lissages de spikes puis des reclassements de courbes, partie 1

```

464 #exponentielle]
465 r2e3 <- numeric(1)
466 par(mfrow=c(1,1))
467 for (i in seq2){
468
469   mois <- 1:35
470   v <- as.vector(puits[,i])
471
472   col="black"
473   if (v[36] == "Good"){
474     col = "red"
475   }else if (v[36] == "medium"){
476     col = "green"
477   }else
478     col = "blue"
479
480   v <- as.numeric(v[1:35])
481
482   #on plot pas les 0 qui sont des ND (d'apres moi)
483   nd <- which(v %in% 0)
484   v <- v[v != 0]
485   mois <- mois[!mois %in% nd]
486
487   smooth <- loess(v~mois)
488
489   nd2 <- which(smooth$fitted <= 0)
490   smooth$fitted <- smooth$fitted[smooth$fitted > 0]
491   mois <- mois[!mois %in% nd2]
492
493   expfit <- lm(log(smooth$fitted) ~ mois)
494
495   #tracé de la figure 1 : les données de production
496   if (i==1){
497     plot(mois,exp(predict(expfit)),type="l",col=col, ylab="gas prod", main="Régression exponentielle avec smooth",y
498   } else {
499     lines(mois,exp(predict(expfit)),type="l",col=col)
500   }
501
502   rsquared = summary(expfit)$adj.r.squared
503   print (rsquared)
504   r2e3=r2e3+summary(expfit)$adj.r.squared
505 }
506

```

FIGURE 22 – Code de la question 5 des lissages de spikes puis des reclassements de courbes, partie 2