

Projet d'Apprentissage Supervisé

Wafa AYAD - Sina AHMADI

le 16 janvier 2017

Abstract

Le projet actuel représente les travaux des auteurs en tant que compte rendu du projet demandé par M. Lazhar LABIOD. Cela correspond au module d'Apprentissage Supervisé de Master2 Intelligence Artificielle à l'Université Paris Descartes.

Après avoir traité les jeux de données, certains algorithmes d'apprentissage supervisé sont appliqués sur les données. Ce document résume la partie théorique du travail résultant de la partie pratique jointe au même dossier rendu.

Table des matières

Introduction	2
Partie I : Étude sur données synthétiques	2
1 Analyse : <i>Aggregation.txt</i>, <i>Flame.txt</i> et <i>Spiral.txt</i>	2
1.1 <i>Flame.txt</i>	2
1.2 <i>Spiral.txt</i>	3
1.3 <i>Aggregation.txt</i>	4
2 Comparaison des modèles de classification supervisée	4
2.1 Apprentissage sur <i>Flame.txt</i>	4
2.2 Apprentissage sur <i>Spiral.txt</i>	5
2.3 Apprentissage sur <i>Aggregation.txt</i>	5
3 Synthèse et conclusion	5
Partie II: Quelques points importants	5
1. Étude exploratoire préliminaire	5
2. Différents algorithmes, différents résultats	5
3. Une courbe explicative	6
4 Annexe	6
Références	12

Introduction

Différentes approches d'apprentissage automatique supervisé existent. Le type de données, leur volume et la disponibilité suffisante de leurs caractéristiques sont des facteurs centraux du choix des algorithmes performants permettant d'avoir un bon rendu d'apprentissage, lequel est nécessaire pour l'étape de prédiction par la suite.

Dans ce rapport, nous allons synthétiser une étude comparative entre plusieurs méthodes d'apprentissage sur des données synthétiques dont les labels sont connus. Ensuite, nous appliquons les modèles de classements adéquats sur un cas des liens d'une banque dans le but d'estimer un score d'appétence à la carte *VISA Premier*.

Ce document est constitué de deux étapes majeurs, la première est consacrée à l'étude comparative entre les algorithmes d'apprentissage supervisé appliqués sur trois jeux de données fournis (*flame.r*, *spiral.r* et *aggregation.r*). Quant à la deuxième partie, et après une étude exploratoire préliminaire de *VisaPremier.txt*, nous discutons l'application des approches que nous jugeons convenables aux données réelles fournies.

Partie I : Étude sur données synthétiques

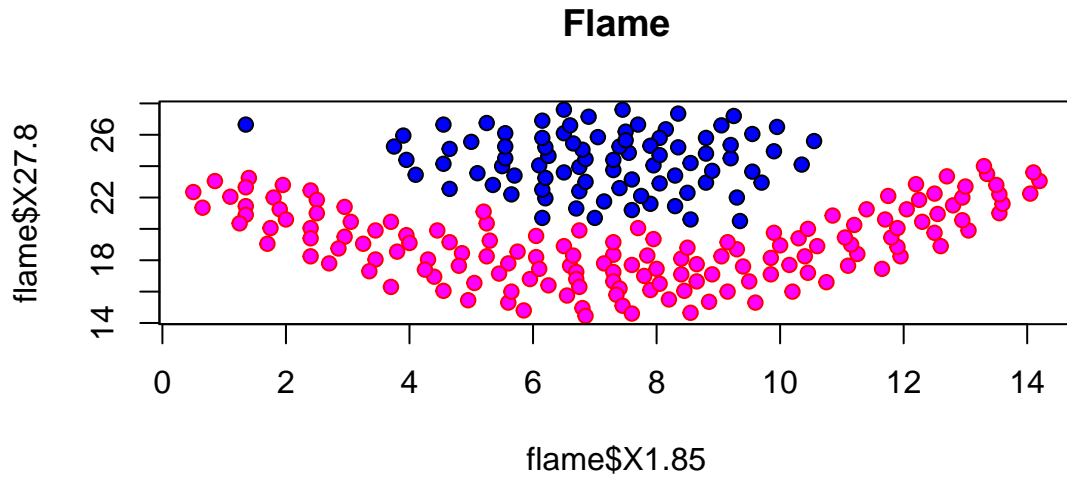
Dans cette partie nous réalisons une étude comparative des différentes approches de classification supervisée vu en cours telles que :

- Logistic Regression
- Linear Discriminant Analysis (LDA)
- Quadratic Discriminant Analysis (QDA)
- Naive Bayes
- K-nearest Neighbours (KNN)
- Support Vector Machine (SVM)

1 Analyse : *Aggregation.txt*, *Flame.txt* et *Spiral.txt*

1.1 *Flame.txt*

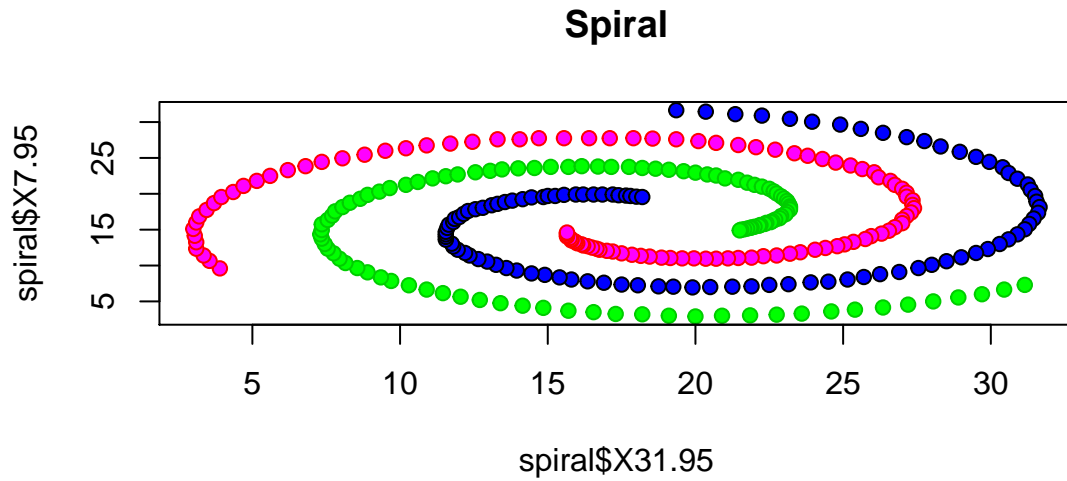
Ces données sont représentées par une matrice de 239 individus décrits par 2 variables de valeurs quantitatives continues. Ces individus sont divisés en 2 classes. Il n'y a pas des valeurs manquantes. Les données n'ont pas besoin de méthodes de prétraitement et peuvent être utilisées directement pour l'apprentissage et la prédiction.



Comme est illustré dans la figure précédente, les données forment 2 classes plus au moins séparées. Une ligne sous forme d'une parabole peut les bien séparer. Les algorithmes linéaires ne pourront pas trouver une bonne marge entre les deux classes dans ce cas-là.

1.2 *Spiral.txt*

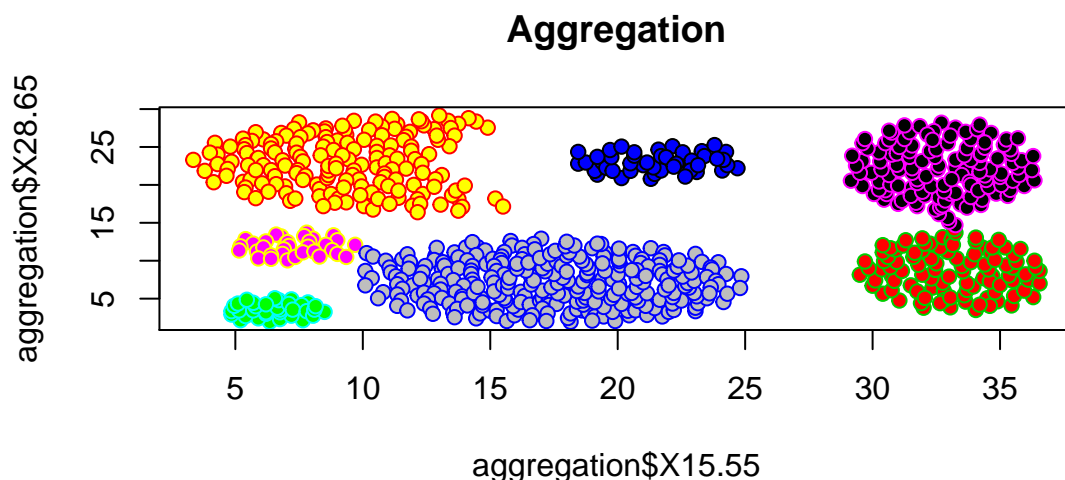
Les 311 individus de cette base sont bivariés, et sont divisés en 3 classes. Aucune des méthodes de prétraitement n'est nécessaire à appliquer sur ces données.



La visualisation des données du fichier *spiral*, comme le montre la figure au dessus, nous permet de distinguer 3 classes spirales bien séparées, les classificateurs linéaires ne pourront pas les distinguer. Cependant, d'autres méthodes non linéaires devraient pouvoir trouver ces 3 classes.

1.3 Aggregation.txt

Cette base de données contient 788 individus représentés par 2 variables dont les valeurs sont quantitatives continues. Nous remarquons qu'il n'y a pas de valeurs manquantes. Les intervalles des valeurs pour les deux variables sont comparables, donc nous n'avons pas besoin de prétraitement tels que la mise en échelle ni la normalisation pour entamer la phase d'apprentissage. Les éléments de cette base sont regroupés en 7 classes différentes comme est illustré dans la figure suivante.



Le schéma ci-dessus montre une bonne séparation des 7 classes, sauf quelques rares points qui sont à la limite d'une classe et une autre. Il est clair que les séparateurs linéaires (tel que LDA) devraient pouvoir détecter facilement ces 7 classes.

2 Comparaison des modèles de classification supervisée

Dans cette section, nous appliquons les différents modèles de classement des données sur les 3 échantillons fournis. Ensuite, nous synthétisons une étude comparative entre les performances de chaque modèle sur chaque base de données dans trois tableaux. Le code utilisé est disponible dans les fichiers joints *flame.r*, *spiral.r* et *aggregation.r*.

Pour appliquer les différentes approches citées auparavant, nous avons consacré 70% des données de chaque base pour la phase d'apprentissage et 30% pour le test.

2.1 Apprentissage sur *Flame.txt*

Le tableau ci-après résume la précision obtenue par les méthodes citées auparavant sur les données *Flame*. Pour le KNN, nous avons utilisé une boucle (de 1 à 10) pour trouver empiriquement le bon K qui donne par la suite de bonnes performances. Quant au SVM, nous l'avons testé avec plusieurs Kernels à savoir: *linear*, *polynomial*, *radial* et *sigmoid*. Les précisions trouvées varient entre celles qui sont bonnes, moyennes et moins bonnes.

Méthode	Logistic Regression	LDA	QDA	Naive Bayes	KNN (k=4)	SVM(Radial Guassien)
Précision	0.46	0.36	0.56	0.67	0.92	0.99

Pour une bonne classification de la base *Flame*, il est recommandable d'utiliser un SVM avec un noyau radial gaussien car sur un échantillon de 100 individus il ne se trompe que dans le classement d'un seul élément. Le 4NN a aussi un moins d'erreurs de classification. Par contre, LDA, QDA et la régression logistique ne donnent pas de bonnes performances dans l'apprentissage sur les données *Flame*.

2.2 Apprentissage sur *Spiral.txt*

Méthode	Logistic Regression	LDA	QDA	Naive Bayes	KNN (k=4)	SVM(Radial Guassien)
Précision	0.33	0.21	0.0	0.0	0.21	0.87

2.3 Apprentissage sur *Aggregation.txt*

Méthode	Logistic Regression	LDA	QDA	Naive Bayes	KNN (k=4)	SVM(Radial Guassien)
Précision	0.06	0.36	0.12	0.67	0.92	0.99

3 Synthèse et conclusion

L'exécution des méthodes d'apprentissage citées auparavant a donné des résultats différents d'un modèle à un autre, cela est dû au fait que chaque algorithme a ses points forts et points faibles, car le type, la dispersion et le volume des données sont les facteurs sur lesquels les modèles s'appuient.

Partie II: Quelques points importants

1. Étude exploratoire préliminaire

Cette étape joue un rôle primordial pour le reste du projet. Les jeux de données ayant des colonnes et des individus, ne garantissent pas forcément la qualité et la cohérence entre eux. Selon nous c'était exigeant de pré-traiter les jeux de données initiaux pour ce projet.

On observe certains phénomènes en analysant des données. Tout d'abord, les données manquantes se trouvaient dans plusieurs colonnes. Les données manquantes sont les données qui fournissent pas assez de données pour qu'elles soient considérées pertinentes. Les colonnes ayant des valeurs non cohérentes, répétées et erronées peuvent aussi être supprimées. Le raffinement de la bruyance est aussi les variables peu-représentatives peuvent aussi mener à avoir un jeu de données plus fiable.

Initialement il existait 47 variables et après la phase du pré-traitement, on a pu extraire 34 variables, en plus la variable de la classe.

La normalisation des données était aussi une des tâches initiales qu'on a effectuée sur le jeu des données. Ayant $\sum_{i,j} \frac{x_{ij} - \bar{x}_i}{\sigma^2}$, on pouvait centrer les données autour de la moyenne des individus et ensuite les normaliser par l'écart-type approprié.

2. Différents algorithmes, différents résultats

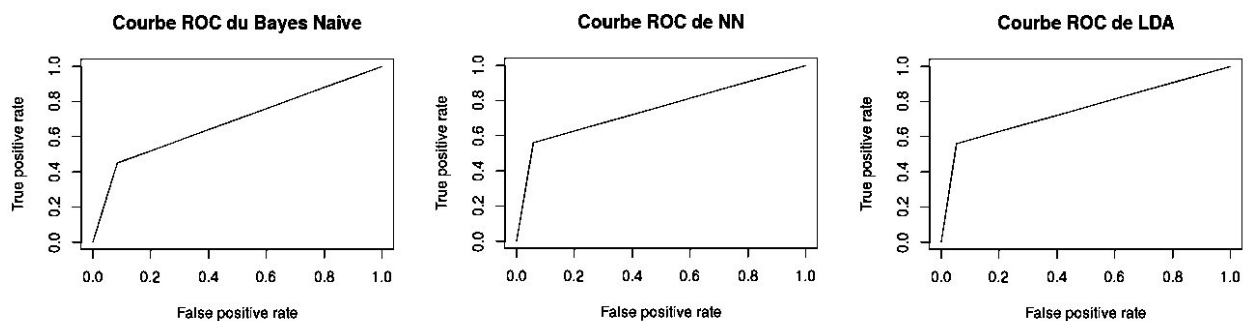
Comme on a vu dans la première partie, chacun des algorithmes peut mener aux résultats qui ne seront pas forcément identiques. Alors que l'algorithme de LDA pour les trois classes, propose une valeur d'estimation

basse mais pertinente, l'algorithme de Naive Bayes ne montre aucune qualité pour la classe Spiral. Pour la même classe on voit que l'algorithme de QDA est aussi égal à 0. Ces valeurs indiquent le lien entre l'essence des données et l'application de chacun des algorithmes à ses propres façons.

Étant les résultats obtenus, on peut bien avouer que les résultats d'algorithme de SVM représente le meilleur. En général, l'algorithme de SVM est un algorithme performant dans le domaine d'analyse de données.

3. Une courbe explicative

Afin de mieux comparer les résultats obtenus par les différents algorithmes, on a essayé de schématiser les résultats de chacun par le diagramme suivant. Le diagramme de ROC (Receiver Operating Characteristic) est un diagramme pour illustrer le rapport sensibilité/spécificité afin d'évaluer la performance de la classification. Les résultats des algorithmes Naive Bayes, LDA et KNN sont représentés dans les diagrammes suivants.



En comparaison des courbes de ROC, on aperçoit que les deux courbes en KNN et LDA partagent les mêmes informations. Plus la courbe est élevée, mieux c'est le modèle. Ce qu'on observe témoigne le fait que les deux algorithmes KNN et LDA ont une qualité meilleure que celle de Naive Bayes.

4 Annexe

L'intégralité des codes en R est disponible dans le dossier contenant ce document. Les 3 parties des codes qui viennent ci-dessous sont chacune pour agrégation, flamme et spiral. On a essayé de bien mettre les commentaires afin d'avoir une meilleure lisibilité pendant l'exécution.

Le contenu du fichier aggregation.r :

```

1
2 # Redirection vers le bon path
3 setwd("F:/M2-MLDS/Apprentissage supervisé/Projet_app_sup_MLDS16")
4
5 #Chargement des données
6 aggregation <- read.table("Aggregation.txt", header = T, sep="\t")
7
8
9 #Aggregation
10 dim(aggregation)
11 plot(aggregation$X15.55, aggregation$X28.65, col= aggregation$X2,
12      pch=21, bg=c("blue", "yellow", "red", "gray", "green", "black", "magenta")
13      [as.numeric(aggregation$X2)], main="Aggregation")
14
15 dim(aggregation)
16 #aggregation[, -3] : Les variables de aggregation
17 #aggregation[, 3] : les classes (2 classes) de aggregation
18
19 #Nous éclatons aggregation en 2 échantillons
20 #Apprentissage : 89% aggregation ==> 700 ind pour le training set
21 #Test : 11% aggregation ==> 87 pour test set
22
23 aggregation_train = aggregation[1:87,]
24 aggregation_test = aggregation[88:787,]

```

```

25
26 #Vérification qu'on a bien 61 pour le test et 250 pour l'apprentissage
27 dim(agggregation_train)
28 dim(agggregation_test)
29
30 #####
31 # Apprentissage & test #
32 #####
33
34
35 #####Regression Logistique#####
36 #Apprentissage sur l'ensemble de aggregation_train
37 aggregation.logReg<- glm(agggregation_train[,3] ~.,data=agggregation_train[,,-3])
38 #Prédiction
39 aggregation.logReg.pred <- predict(agggregation.logReg,newdata=agggregation_test[,,-3],type='response')
40 #Selon le seuil,
41 aggregation.logReg.pred <- ifelse(agggregation.logReg.pred > 0.5,1,0)
42 #Calcul de la précision
43 aggregation.logReg.erreurClass = mean(agggregation.logReg.pred != aggregation_test[,3])
44 aggregation.logReg.accuracy = 1 - aggregation.logReg.erreurClass
45 print(paste("Précision Logistic Regression =",
46 round(agggregation.logReg.accuracy, digits = 2)*100,"%"))
47
48 #La regression logistique donne une précision de 6% :(
49
50 #####LDA#####
51 library(MASS)
52 #Construction de discrimination linéaire
53 aggregation.lda <- lda(agggregation_train[,,-3], agggregation_train[,3] )
54 #Prediction selon le modèle linéaire
55 aggregation.lda.pred <- predict(agggregation.lda,agggregation_test[,,-3])
56 #Table de confusion
57 table(agggregation.lda.pred$class,agggregation_test[,3])
58
59 #Calcul de la précision que la LDA donne
60 aggregation.lda.erreurClass <- mean(agggregation.lda.pred$class != aggregation_test[,3])
61 aggregation.lda.accuracy = 1 - aggregation.lda.erreurClass
62 print(paste("Précision LDA =",round(agggregation.lda.accuracy, digits = 2)*100,"%"))
63 #LDA donne une précision de 36% :(
64
65 #####KNN#####
66 library(class)
67 #Pour trouver de meilleure performance, nous devrions tester l'algorithme KNN pour
68 #différentes valeurs de K, nous ne sauvegardons que le k donnant la meilleure précision.
69 k.optim = 0
70 aggregation.KNN.accuracy = 0
71 for (i in 1:10) {
72 aggregation.KNN = knn(agggregation_train[,,-3], agggregation_test[,,-3], cl = aggregation_train[,3] , k = i
73 )
74 tmp = 1 - ( sum(agggregation.KNN != aggregation_test[,3]) / length(agggregation_test[,3]) )
75 if(tmp > aggregation.KNN.accuracy){
76 aggregation.KNN.accuracy = tmp
77 k.optim = i
78 }
79 }
80 print(paste("Précision ",k.optim ,"NN =",round(agggregation.KNN.accuracy, digits = 2)*100,"%"))
81 #La meilleur valeur de K (allant de 1:10 ou même 1:20) est 4
82 #La précision de 4NN est 92% meilleure que celle de la LDA; :)
83
84
85 #####Classificateur Baysien Naif#####
86
87 library(e1071)
88 # apprentissage sur le training set
89 aggregation.NaiveB = naiveBayes(as.factor(agggregation_train[,3]) ~ ., data = aggregation_train[,,-3])
90 #Prediction sur l'échantillon test
91 aggregation.NaiveB.pred = predict(agggregation.NaiveB, agggregation_test[,,-3])
92 #Table de confusion
93 table(agggregation.NaiveB.pred,agggregation_test[,3])
94 #Calcul de la précision
95 aggregation.NaiveB.erreurClass <- mean(agggregation.NaiveB.pred != aggregation_test[,3])
96 aggregation.NaiveB.accuracy = 1 - aggregation.NaiveB.erreurClass
97 print(paste("Précision Naive Bayes =",round(agggregation.NaiveB.accuracy, digits = 2)*100,"%"))
98 # Le score obtenu est de 67% mieux que la LDA mais moins bon que KNN :/
99
100
101 #####QDA#####
102 #Construction de discrimination quadratique
103 aggregation.qda <- qda(agggregation_train[,,-3], agggregation_train[,3])
104 #Prediction selon le modèle quadratique

```

```

105 aggregation.qda.pred <- predict(aggregation.qda,aggregation_test[,-3])
106 #Table de confusion
107 table(aggregation.qda.pred$class, aggregation_test[,3])
108 #Calcul de la précision
109 aggregation.qda.erreurClass <- mean(aggregation.qda.pred$class != aggregation_test[,3])
110 aggregation.qda.accuracy = 1 - aggregation.qda.erreurClass
111 print(paste("Précision QDA =",round(aggregation.qda.accuracy, digits = 2)*100,"%"))
112
113 #Contrairement é ce qu'on aurait pu penser la QDA ne donne pas de très bon résultats
114 #Seulement 12% de précision (mieux que LDA mais moins bonne que KNN et Naive Bayes :)
115
116
117 #####SVM#####
118 #Définition des kernels nécessaires au SVM
119 kernels = c("linear","polynomial","radial","sigmoid")
120 #La fonction tune estime les paramètres
121 aggregation.svmTune = tune(svm, train.x = aggregation_train[,-3], train.y = aggregation_train[,3],
122                           validation.x = aggregation_test[,-3], validation.y = aggregation_test[,3],
123                           ranges = list(kernel = kernels))
124
125 # précision du meilleur modèle
126 aggregation.svm.accuracy = 1 - aggregation.svmTune$best.performance
127 print(paste("Précision SVM (",aggregation.svmTune$best.parameters$kernel,") =",
128           round(aggregation.svm.accuracy, digits = 2)*100,"%"))
129 #Avec le kernel gaussien radial, nous obtenons un taux de précision supérieur é 99% :D

```

Le contenu du fichier spiral.r :

```

1 # Redirection vers le bon path
2 setwd("F:/M2-MLDS/Apprentissage supervisé/Projet_app_sup_MLDS16")
3
4 #Chargement des données
5 spiral <- read.table("spiral.txt",header = T, sep = "\t", fill = T)
6
7 #spiral
8 dim(spiral)
9 plot(spiral$X31.95, spiral$X7.95, col= spiral$X3, pch=21,bg=c("blue", "magenta", "green")[as.numeric(
   spiral$X3)], main="spiral")
10
11 #spiral[,,-3] : Les variables de spiral
12 #spiral[,3] : les classes (2 classes) de spiral
13
14 #Nous éclatons spiral en 2 échantillons (é 1/3 de la classe 1, 1/3 de la classe 2 et 1/3 de la classe 3)
15 #Apprentissage : 80% spiral ==> 250 ind pour le training set
16 #Test : 20% spiral ==> 61 pour test set
17
18 spiral_train = rbind(spiral[21:105,],spiral[126:206,], spiral[228:311, ])
19 spiral_test = rbind(spiral[1:20,], spiral[106:125,], spiral[207:227, ])
20
21 #Vérification qu'on a bien 61 pour le test et 250 pour l'apprentissage
22 dim(spiral_train)
23 dim(spiral_test)
24
25 #La précision, et la spécificité semblent étre des bonnes mesures pour ésitmer les
26 #méthodes sur ce jeu de données
27 #Le nombre d'individu par classe est assez similaire
28
29 #####
30 # Apprentissage & test #
31 #####
32
33
34 #####Regression Logistique#####
35 #Apprentissage sur l'ensemble de spiral_train
36 spiral.logReg<- glm(spiral_train[,3] ~.,data=spiral_train[,-3])
37 #Prédiction
38 spiral.logReg.pred <- predict(spiral.logReg,newdata=spiral_test[,-3],type='response')
39 #Selon le seuil,
40 spiral.logReg.pred <- ifelse(spiral.logReg.pred > 0.5,1,0)
41 #Calcul de la précision
42 spiral.logReg.erreurClass = mean(spiral.logReg.pred != spiral_test[,3])
43 spiral.logReg.accuracy = 1 - spiral.logReg.erreurClass
44 print(paste("Précision Logistic Regression =",
45           round(spiral.logReg.accuracy, digits = 2)*100,"%"))
46
47 #La regression logistique donne une précision de 33% :(
48
49 #####LDA#####
50 library(MASS)
51 #Construction de discrimination linéaire

```



```

52 spiral.lda <- lda(spiral_train[,-3], spiral_train[,3] )
53 #Prediction selon le modèle linéaire
54 spiral.lda.pred <- predict(spiral.lda,spiral_test[,-3])
55 #Table de confusion
56 table(spiral.lda.pred$class,spiral_test[,3])
57
58 #Calcul de la précision que la LDA donne
59 spiral.lda.erreurClass <- mean(spiral.lda.pred$class != spiral_test[,3])
60 spiral.lda.accuracy = 1 - spiral.lda.erreurClass
61 print(paste("Précision LDA =",round(spiral.lda.accuracy, digits = 2)*100,"%"))
62 #LDA donne une précision de 0% :'(, si t'as dis je calcule pas ça aurait été mieux :/
63
64 #####KNN#####
65 library(class)
66 #Pour trouver de meilleure performance, nous devrions tester l'algorithme KNN pour
67 #différentes valeurs de K, nous ne sauvegardons que le k donnant la meilleure précision.
68 k.optim = 0
69 spiral.KNN.accuracy = 0
70 for (i in 1:10) {
71   spiral.KNN = knn(spiral_train[,-3], spiral_test[,3], cl = spiral_train[,3] , k = i)
72   tmp = 1 - ( sum(spiral.KNN != spiral_test[,3]) / length(spiral_test[,3]) )
73   if(tmp > spiral.KNN.accuracy){
74     spiral.KNN.accuracy = tmp
75     k.optim = i
76   }
77 }
78
79 print(paste("Précision ",k.optim ,"NN =",round(spiral.KNN.accuracy, digits = 2)*100,"%"))
80 #La meilleur valeur de K (allant de 1:10 ou même 1:20) est 1, ce qui est logique car
81 #les spirales sont loin l'une des autres !
82 #La précision de 1NN est 21% , Deéue :( !
83
84
85 #####Classificateur Baysien Naif#####
86
87 library(e1071)
88 # apprentissage sur le training set
89 spiral.NaiveB = naiveBayes(as.factor(spiral_train[,3]) ~ ., data = spiral_train[,-3])
90 #Prediction sur l'échantillon test
91 spiral.NaiveB.pred = predict(spiral.NaiveB, spiral_test[,-3])
92 #Table de confusion
93 table(spiral.NaiveB.pred,spiral_test[,3])
94 #Calcul de la précision
95 spiral.NaiveB.erreurClass <- mean(spiral.NaiveB.pred != spiral_test[,3])
96 spiral.NaiveB.accuracy = 1 - spiral.NaiveB.erreurClass
97 print(paste("Précision Naive Bayes =",round(spiral.NaiveB.accuracy, digits = 2)*100,"%"))
98 # Le score obtenu est de 0% (au moins toi t'es honête -_-)
99
100
101 #####QDA#####
102 #Construction de discrimination quadratique
103 spiral.qda <- qda(spiral_train[,-3], spiral_train[,3])
104 #Prediction selon le modèle quadratique
105 spiral.qda.pred <- predict(spiral.qda,spiral_test[,-3])
106 #Table de confusion
107 table(spiral.qda.pred$class, spiral_test[,3])
108 #Calcul de la précision
109 spiral.qda.erreurClass <- mean(spiral.qda.pred$class != spiral_test[,3])
110 spiral.qda.accuracy = 1 - spiral.qda.erreurClass
111 print(paste("Précision QDA =",round(spiral.qda.accuracy, digits = 2)*100,"%"))
112
113 # 0% :/
114
115
116 #####SVM#####
117 #Définition des kernls nécessaires au SVM
118 kernels = c("linear","polynomial","radial","sigmoid")
119 #La fonction tune estime les paramètres
120 spiral.svmTune = tune(svm, train.x = spiral_train[,-3], train.y = spiral_train[,3],
121                      validation.x = spiral_test[,-3], validation.y = spiral_test[,3],
122                      ranges = list(kernel = kernels))
123
124 # précision du meilleur modèle
125 spiral.svm.accuracy = 1 - spiral.svmTune$best.performance
126 print(paste("Précision SVM (",spiral.svmTune$best.parameters$kernel,") =",
127           round(spiral.svm.accuracy, digits = 2)*100,"%"))
128 #Avec le kernel gaussien radial, nous obtenons un taux de précision supérieur é 87% :D

```

Le contenu du fichier flamer.r :

```

1 # Redirection vers le bon path
2 setwd("F:/M2-MLDS/Apprentissage supervisé/Projet_app_sup_MLDS16")
3
4 #Chargement des données
5 flame <- read.table("flame.txt",header = T, sep = "\t", fill = T)
6
7 #Dimension de flame
8 dim(flame)
9 #affichage des individus de flame en fonction des 2 variables et en les colorie
10 #en fonction de leur classe
11 plot(flame$X1.85, flame$X27.8, col=flame$X1, pch=21,bg=c("blue", "magenta")
12       [as.numeric(flame$X1)], main="Flame")
13 #Il est clair que les 2 classes sont séparées par une courbe parabolique
14 #(de type  $y=x_0+x_1^2$ )
15
16
17 #flame[, -3] : Les variables de flame
18 #Flame[, 3] : les classes (2 classes) de flame
19
20 #Nous éclatons flame en 2 échantillons (à 1/2 de la classe 1 et 1/2 de la classe 2)
21 #Apprentissage : 84% flame ==> 200 ind pour le training set
22 #Test : 16% flame ==> 39 pour test set
23
24 flame_train = rbind(flame[1:133,], flame[173:240,])
25 flame_test = flame[134:172,]
26
27 #Vérification qu'on a bien 39 pour le test et 200 pour l'apprentissage
28 dim(flame_train)
29 dim(flame_test)
30
31 #La précision, et la spécificité semblent être des bonnes mesures pour évaluer les
32 #méthodes sur ce jeu de données
33 #Le nombre d'individu par classe est assez similaire
34
35 #####
36 # Apprentissage & test
37 #####
38
39
40 #####Regression Logistique#####
41 #Apprentissage sur l'ensemble de flame_train
42 flame.logReg <- glm(flame_train[, -3] ~., data=flame_train[, -3])
43 #Prédiction
44 flame.logReg.pred <- predict(flame.logReg, newdata=flame_test[, -3], type='response')
45 #Selon le seuil,
46 flame.logReg.pred <- ifelse(flame.logReg.pred > 0.5, 1, 0)
47 #Calcul de la précision
48 flame.logReg.erreurClass = mean(flame.logReg.pred != flame_test[, 3])
49 flame.logReg.accuracy = 1 - flame.logReg.erreurClass
50 print(paste("Précision Logistic Regression =",
51            round(flame.logReg.accuracy, digits = 2)*100, "%"))
52
53 #La regression logistique donne une précision de 46% :(
54
55 #####LDA#####
56 library(MASS)
57 #Construction de discrimination linéaire
58 flame.lda <- lda(flame_train[, -3], flame_train[, 3])
59 #Prediction selon le modèle linéaire
60 flame.lda.pred <- predict(flame.lda, flame_test[, -3])
61 #Table de confusion
62 table(flame.lda.pred$class, flame_test[, 3])
63
64 #Calcul de la précision que la LDA donne
65 flame.lda.erreurClass <- mean(flame.lda.pred$class != flame_test[, 3])
66 flame.lda.accuracy = 1 - flame.lda.erreurClass
67 print(paste("Précision LDA =", round(flame.lda.accuracy, digits = 2)*100, "%"))
68 #LDA donne une précision de 36% :(
69
70 #####KNN#####
71 library(class)
72 #Pour trouver de meilleure performance, nous devrions tester l'algorithme KNN pour
73 #différentes valeurs de K, nous ne sauvegardons que le k donnant la meilleure précision.
74 k.optim = 0
75 flame.KNN.accuracy = 0
76 for (i in 1:10) {
77     flame.KNN = knn(flame_train[, -3], flame_test[, -3], cl = flame_train[, 3], k = i)
78     tmp = 1 - (sum(flame.KNN != flame_test[, 3]) / length(flame_test[, 3]))
79     if(tmp > flame.KNN.accuracy){
80         flame.KNN.accuracy = tmp
81         k.optim = i

```

```

82   }
83 }
84
85 print(paste("Précision ",k.optim,"NN =",round(flame.KNN.accuracy, digits = 2)*100,"%"))
86 #La meilleur valeur de K (allant de 1:10 ou même 1:20) est 4
87 #La précision de 4NN est 92% meilleure que celle de la LDA; :)
88
89
90 #####Classificateur Baysien if#####
91
92 library(e1071)
93 # apprentissage sur le training set
94 flame.NaiveB = naiveBayes(as.factor(flame_train[,3]) ~ ., data = flame_train[, -3])
95 #Prediction sur l'échantillon test
96 flame.NaiveB.pred = predict(flame.NaiveB, flame_test[, -3])
97 #Table de confusion
98 table(flame.NaiveB.pred, flame_test[,3])
99 #Calcul de la précision
100 flame.NaiveB.erreurClass <- mean(flame.NaiveB.pred != flame_test[,3])
101 flame.NaiveB.accuracy = 1 - flame.NaiveB.erreurClass
102 print(paste("Précision Naive Bayes =",round(flame.NaiveB.accuracy, digits = 2)*100,"%"))
103 # Le score obtenu est de 67% mieux que la LDA mais moins bon que KNN :/
104
105
106 #####QDA#####
107 #Construction de discrimination quadratique
108 flame.qda <- qda(flame_train[, -3], flame_train[,3])
109 #Prediction selon le modèle quadratique
110 flame.qda.pred <- predict(flame.qda, flame_test[, -3])
111 #Table de confusion
112 table(flame.qda.pred$class, flame_test[,3])
113 #Calcul de la précision
114 flame.qda.erreurClass <- mean(flame.qda.pred$class != flame_test[,3])
115 flame.qda.accuracy = 1 - flame.qda.erreurClass
116 print(paste("Précision QDA =",round(flame.qda.accuracy, digits = 2)*100,"%"))
117
118 #Contrairement à ce qu'on aurait pu penser la QDA ne donne pas de très bon résultats
119 #Seulement 56% de précision (mieux que LDA mais moins bonne que KNN et Naive Bayes :/
120
121
122 #####SVM#####
123 #Définition des kernls nécessaires au SVM
124 kernels = c("linear","polynomial","radial","sigmoid")
125 #La fonction tune estime les paramètres
126 flame.svmTune = tune(svm, train.x = flame_train[, -3], train.y = flame_train[,3],
127                     validation.x = flame_test[, -3], validation.y = flame_test[,3],
128                     ranges = list(kernel = kernels))
129
130 # précision du meilleur modèle
131 flame.svm.accuracy = 1 - flame.svmTune$best.performance
132 print(paste("Précision SVM (",flame.svmTune$best.parameters$kernel,") =",
133           round(flame.svm.accuracy, digits = 2)*100,"%"))
134 #Avec le kernel gaussien radial, nous obtenons un taux de précision supérieur é 99% :D

```

References

- [1] Keil, P., Belmaker, J., Wilson, A. M., Unitt, P., & Jetz, W. (2013). Downscaling of species distribution models: a hierarchical approach. *Methods in Ecology and Evolution*, 4(1), 82-94.
- [2] Team, R. Core. "R: A language and environment for statistical computing." (2013): 409.
- [3] Halpern, B. S., Regan, H. M., Possingham, H. P., & McCarthy, M. A. (2006). Accounting for uncertainty in marine reserve design. *Ecology Letters*, 9(1), 2-11.

.