

TP Analyse de données - Apprentissage supervisé (Classification)

INTRODUCTION

Nous allons utiliser le logiciel R ([documentation](#) ; possibilité d'utiliser Rstudio). Dans un dossier *AnalyseDeDonnees* créer deux sous dossiers :

- Code - dans lequel vous placerez vos fichiers de code
- Data - dans lequel vous placerez les fichiers du répertoire "Data" (à télécharger via Moodle)

Les 3 TP *Analyse de données* présentent différentes méthodes d'analyse de données : le but n'est pas de finir les TP le plus vite possible mais d'analyser les résultats ! Un rapport de 1 page **maximum** vous est demandé pour chacun des 3 TP : ne choisir que les résultats les plus intéressants et les **commenter**.

Ne pas hésiter à utiliser l'aide de R grâce à la commande :

```
1 help(...)
```

CLASSIFICATION - APPROCHES DE TYPE PROTOTYPE

0. Télécharger le cours sur la classification (seconde partie du cours).

Méthode des k-plus proches voisins

1. Créer un fichier *classification-kppv.R* et copier les lignes suivantes :

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../TP/TP/Code")
3 # Packages utilisés dans la suite
4 library(class)
5 library(caret)
6 library(ROCR)
7 # Supprimer toutes les variables
8 rm(list=ls(all=TRUE))
9 # Supprimer tous les graphiques déjà présents
10 graphics.off()
```

2. Charger les données (synthétiques) d'**apprentissage** et les séparer en 2 variables : données d'entrée et données de sortie. L'objectif est de déterminer si y vaut 1 ou 2 connaissant x_1 et x_2 .

```
1 # Lecture des données d'apprentissage
2 data_train <- read.table("../Data/synth_train.txt", header=T, sep="\t");
3 print(data_train)
4 # Séparation des données et de la sortie
5 data_train_x <- data.frame(x1=data_train$x1, x2=data_train$x2)
```

3. Charger les données **test** et les séparer aussi en 2 variables :

```
1 # Lecture des données test
2 data_test <- read.table("../Data/synth_test.txt",header=T,sep="\t");
3 print(data_test)
4 # Séparation des données et de la sortie
5 data_test_x <- data.frame(x1=data_test$x1,x2=data_test$x2)
```

4. Tracer les données (colorées par la sortie y)

```
1 # Graphique des données (colorées par la sortie y)
2 plot(data_train$x1,data_train$x2,col=data_train$y,pch=16)
```

5. Appliquer les k -plus proches voisins sur les données d'apprentissage :

```
1 # nombre de voisins (par ex proche de la racine carré du nombre d'obs)
2 num_of_neigh <- 10
3 data_train_predict <- knn(train=data_train_x,test=data_train_x,
4                           cl=data_train$y,k=num_of_neigh)
5 # Affichage des résultats (étoile)
6 par(new=T)
7 plot(data_train$x1,data_train$x2,col=data_train_predict,pch=8)
8 # Calcul du taux d'erreur
9 error_rate <- mean(data_train_predict != data_train$y)
10 cat("error_rate using train data = ",error_rate)
```

6. Appliquer les k -plus proches voisins sur les données **test** :

```
1 data_test_predict <- knn(train=data_train_x,test=data_test_x,
2                           cl=data_train$y,k=num_of_neigh)
3 # Affichage des données (cercle)
4 plot(data_train$x1,data_train$x2,col=data_train$y,pch=16)
5 par(new=T)
6 # Affichage des résultats (étoile)
7 plot(data_test$x1,data_test$x2,col=data_test_predict,pch=8)
8 # Affichage des vraies valeurs (triangle)
9 par(new=T)
10 plot(data_test$x1,data_test$x2,col=data_test$y,pch=2)
11 # Calcul du taux d'erreur
12 error_rate <- mean(data_test_predict != data_test$y)
13 cat("error_rate using test data = ",error_rate)
```

7. Afficher la matrice de confusion (définir les vrais pos./nég. et les faux pos./nég.). Imaginons que $y = 1$ (resp. $y = 2$) signifie que le patient est malade (resp. sain)? Est-ce une bonne méthode? Est-ce qu'il vaut mieux dans ce cas là avoir des faux positifs ou des faux négatifs?

```
1 # Matrice de confusion
2 confmat = table(data_test_predict,data_test$y)
3 print("Confusion Matrix")
4 print(confmat)
5 # vrais positifs + vrais négatifs + faux positifs + faux négatifs
6 TP = confmat[1,1]; TN = confmat[2,2]; FP = confmat[1,2]; FN = confmat[2,1];
```

8. Définir la sensibilité, la spécificité, la précision et la prévalence. Les étudier en considérant l'exemple de la question précédente.

```
1 # Sensibilité (sensitivity ; TPR = true positive rate)
2 TPR = TP/(TP+FN)
3 cat("TPR",TPR,"\n")TP
4 # Spécificité (specificity ; TNR = true negative rate)
5 TNR = TN/(TN+FP)
6 cat("TNR",TNR,"\n")
7 # Précision (precision ; positive predictive value)
8 PPV = TP/(TP+FP)
9 cat("PPV",PPV,"\n")
10 # se compare à la prévalence (prevalence)
11 cat("Prev =",length(data_test$y[data_test$y==1])/length(data_test$y),"\n")
```

9. Calculer la précision avec la F-measure (dit F-score).

```
1 cat("F-score = ",2 * TPR * PPV / (TPR+PPV),"\n")
```

10. Nous allons maintenant étudier sur ce petit exemple l'importance du nombre de voisins. Pour cela, nous allons représenter graphiquement la frontière de décision pour $k = 1, 5, 10, 15, 20, 30$ voisins. Commencer par construire une grille de points.

```
1 # Explication visuelle de l'importance de la valeur de k (nb de voisins)
2 # Construction de la grille
3 gridx1 <- seq(from=min(data_train$x1),to=max(data_train$x1),length.out=50)
4 gridx2 <- seq(from=min(data_train$x2),to=max(data_train$x2),length.out=50)
5 grid <- expand.grid(x1 = gridx1, x2 = gridx2)
6 data_grid_x <- data.frame(x1=grid[,1],x2=grid[,2])
```

11. Appliquer la méthode des k-plus proches voisins pour $k = 1, 5, 10, 15, 20, 30$ sur la grille. Que pensez vous des résultats? Que se passe t-il au niveau de la frontière? Réfléchir à un moyen de trouver comment choisir k (en utilisant que les données d'apprentissage!).

```
1 # k plus proches voisins avec application sur les données de la grille
2 num_of_neigh_grid <- c(1,5,10,15,20,30)
3 par(mfrow=c(2,length(num_of_neigh_grid)/2))
4 for (i in 1:length(num_of_neigh_grid))
5 {
6   num_of_n <- num_of_neigh_grid[i]
7   data_g_pr <- knn(train=data_train_x,test=data_grid_x,
8                   cl=data_train$y,k=num_of_n)
9   plot(data_train$x1,data_train$x2,col=data_train$y,pch=16)
10  title(paste("num of neighbours = ",toString(num_of_n)))
11  par(new=T)
12  plot(data_grid_x$x1,data_grid_x$x2,col=data_g_pr,pch=8,cex=0.5,ann=FALSE)
13 }
```

12. Avec la méthode des k-plus proches voisins il est très facile de déterminer la probabilité associée au résultat. En effet supposons un nouveau point dont les valeurs de y pour les 7-plus proches voisins sont :

$$y = 1, y = 2, y = 2, y = 2, y = 1, y = 2, y = 1.$$

Comme il y a 4 voisins sur 7 avec $y = 2$, on va classer ce nouveau point à $y = 2$ mais on peut aussi très facilement déterminer la probabilité de bon classement : $4/7 \sim 0.6$.

```
1 # k plus proches voisins avec les probas
2 data_test_predict_with_proba <- knn(train=data_train_x,test=data_test_x,
3                                     cl=data_train$y,k=num_of_neigh,prob=TRUE)
```

13. À partir de cette probabilité p et du résultat ($y = 1$ ou $y = 2$), il est possible de déterminer un score, c'est-à-dire de déterminer quelle est la probabilité que le patient soit malade (c'est-à-dire $y = 1$). Si $y = 1$ (resp. $y = 2$) on va considérer p (resp. $1 - p$).

```
1 # Calcul du score
2 score <- attr(data_test_predict_with_proba, "prob")
3 score <- ifelse(data_test_predict_with_proba == "1", 1-score, score)
```

14. À partir de ce score il est possible de tracer une courbe ROC. La droite $y = x$ est aussi tracée : celle-ci correspond au classifieur aléatoire.

```
1 pred_knn <- prediction(score, data_test$y)
2 perf <- performance(pred_knn, "tpr", "fpr")
3 plot(perf,colorize=TRUE)
4 par(new=T)
5 plot(c(0,1),c(0,1),type="l",ann=FALSE)
```

15. Pour connaître la performance de ce classifieur, il est possible de calculer l'aire sous la courbe.

```
1 # Aire sous la courbe
2 AUC <- performance(pred_knn, "auc")@y.values[[1]]
3 cat("AUC = ", AUC)
```

16. Il est possible de déterminer le meilleur seuil. Le repérer aussi sur la courbe ROC.

```
1 # Choix du seuil
2 result <- NULL
3 threshold <- seq(0,1,len=11)
4 for (s in threshold)
5 {
6   test <- as.integer(score>=s)+1
7   result <- c(result,1-mean(test != data_test$y))
8 }
9 plot(threshold,result,type="l")
10 cat("Meilleur seuil ", threshold[which.max(result)])
```

Arbres de classification

1. Créer un fichier *classification-CART.R* et copier les lignes suivantes :

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../TP/TP/Code")
3 # Packages utilisés dans la suite
4 library(rpart)
5 # Supprimer toutes les variables
6 rm(list=ls(all=TRUE))
7 # Supprimer tous les graphiques déjà présents
8 graphics.off()
```

2. Les données utilisées sont les mêmes que précédemment.

```
1 # Lecture des données d'apprentissage
2 data_train <- read.table("../Data/synth_train.txt",header=T,sep="\t");
3 plot(data_train$x1, data_train$x2, pch=data_train$y, col=data_train$y,
4      main="Training set")
5 legend("topleft", legend=c("classe1", "classe2"), pch=1:2, col=1:2)
6
7 # Création de la sortie (à mettre sous le format facteur sinon
8 # un modèle de régression est créé)
9 data_train$y <- as.factor(data_train$y)
```

3. Créer l'arbre en utilisant les données d'apprentissage (méthode CART).

```
1 tree <- rpart(y~., data = data_train)
```

4. Afficher les informations sur l'arbre créé. Interpréter les résultats.

```
1 print(tree)
2 summary(tree)
```

5. Tracer l'arbre dans un graphique.

```
1 plot(tree)
2 text(tree)
```

6. Définir la variable correspond à la valeur de séparation sur x_1 .

```
1 # Valeur de séparation de x1
2 split_x1 <- tree$splits[1,4]
```

7. Lire les données test et faire un graphique avec la droite de séparation pour les données d'apprentissage et les données test.

```
1 data_test <- read.table("../Data/synth_test.txt",header=T,sep="\t");
2 par(mfrow=c(1:2))
3 plot(data_train$x1, data_train$x2, pch=data_train$y, col=data_train$y,
4      main="Training set")
5 legend("topleft", legend=c("classe1", "classe2"), pch=1:2, col=1:2)
6 abline(v=split_x1,lty=2,col=4)
7 plot(data_test$x1, data_test$x2, pch=data_test$y, col=data_test$y,
8      main="Test set")
9 legend("topleft", legend=c("classe1", "classe2"), pch=1:2, col=1:2)
10 abline(v=split_x1,lty=2,col=4)
```

8. Évaluer la performance de l'arbre de classification sur les données test.

```
1 # Prediction sur les données test
2 data_test_x <- data.frame(x1=data_test$x1,x2=data_test$x2)
3 tree_predict_data_test <- predict(tree, newdata=data_test_x, type="class")
4
5 # Comparaison des valeurs prédites et des valeurs observées
6 table(tree_predict_data_test, data_test$y)
7
8 # Calcul du taux d'erreur
9 error_rate <- mean(tree_predict_data_test != data_test$y)
10 cat("error_rate using test data = ",error_rate)
```

9. Dans la méthode CART, l'arbre est élagué. Il est possible de construire l'arbre complet en modifiant les valeurs par défaut de la fonction *rpart*.

```
1 # Arbre de longueur maximale
2 tree_max <- rpart(y~., data=data_train, minsplit=2,cp=0)
3 # Tracer l'arbre maximal
4 par(mfrow=c(1,1))
5 plot(tree_max)
6 text(tree_max,use.n = TRUE,all=TRUE)
```

10. En affichant l'arbre complet, interpréter l'élagage précédent.

```
1 # Affichage de l'arbre
2 plotcp(tree_max)
```

Forêts aléatoires

1. Créer un fichier *classification-randomforest.R* et copier les lignes suivantes :

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../TP/TP/Code")
3 # Packages utilisés dans la suite
4 library(randomForest)
5 # Supprimer toutes les variables
6 rm(list=ls(all=TRUE))
7 # Supprimer tous les graphiques déjà présents
8 graphics.off()
```

2. Les données utilisées sont les mêmes que précédemment.

```
1 # Lecture des données d'apprentissage
2 data_train <- read.table("../Data/synth_train.txt",header=T,sep="\t");
3 # Séparation des données et de la sortie !!
4 data_train_x <- data.frame(x1=data_train$x1,x2=data_train$x2)
5 # Création de la sortie (à mettre sous le format facteur sinon
6 # un modèle de régression est créé)
7 data_train_y <- as.factor(data_train$y)
```

3. Appliquer l'algorithme des forêts aléatoires.

```
1 # Forêts aléatoires
2 rf <- randomForest(x = data_train_x, y = data_train_y)
3
4 # Évolution de l'erreur en fonction du nombre d'arbres
5 # Ici ntree est fixé à la valeur par défaut = 500
6 plot(rf$err.rate[,1], type="l")
7
8 # Affichage des résultats
9 print(rf)
```

4. Plusieurs paramètres peuvent être modifiés (ils sont fixés pour le moment à leur valeur par défaut) :

- le nombre d'arbres considérés *ntree* (par défaut à 500)
- la taille de l'échantillon *sampsiz*e (par défaut pas fixé)
- la limite maximale du nombre de noeuds de chaque arbre *maxnodes* (par défaut non fixé)
- quand plus de 2 variables d'entrée, le nombre de variables testé à chaque division *mtry*

Modifier les 3 premiers paramètres pour tester (sur cet exemple très simple les changements ne sont pas majeurs).

5. Afficher l'importance de chaque variable.

```
1 # Importance des variables
2 rf$importance
3 varImpPlot(rf)
```

6. Tester l'algorithme de forêt aléatoire sur les données test. Est-ce mieux que les 2 précédentes méthodes (k-plus proches voisins et CART) ?

```
1 # Lecture des données test
2 data_test <- read.table("../Data/synth_test.txt",header=T,sep="\t");
3 # Séparation des données et de la sortie !!
4 data_test_x <- data.frame(x1=data_test$x1,x2=data_test$x2)
5
6 # Prediction sur les données test
7 rf_predit_data_test <- predict(rf, newdata=data_test_x)
8
9 # Comparaison des valeurs prédites et des valeurs observées
10 table(rf_predit_data_test, data_test$y)
11
12 # Calcul du taux d'erreur
13 error_rate <- mean(rf_predit_data_test != data_test$y)
14 cat("error_rate using test data = ",error_rate)
```

7. Pour comparer les résultats à la méthode des k-plus proches voisins nous allons appliquer les résultats sur une grille pour obtenir la frontière :

```
1 # Résultat de la méthode des forêts aléatoires
2 gridx1 <- seq(from=min(data_train$x1), to=max(data_train$x1), length.out=50)
3 gridx2 <- seq(from=min(data_train$x2), to=max(data_train$x2), length.out=50)
4 grid <- expand.grid(x1 = gridx1, x2 = gridx2)
5 data_grid_x <- data.frame(x1=grid[,1],x2=grid[,2])
6
7 rf_predit_data_grid <- predict(rf, newdata=data_grid_x)
8
9 plot(data_train$x1,data_train$x2,col=data_train$y,pch=16)
10 par(new=T)
11 plot(data_grid_x$x1,data_grid_x$x2,col=rf_predit_data_grid,pch=8,cex=0.5,
12      ann=FALSE)
```


CLASSIFICATION - APPROCHES BASÉES SUR UN MODÈLE

0. Continuer à travailler sur le cours sur la classification (première partie du cours).

Analyse discriminante linéaire et quadratique

1. Créer un fichier *classification-analysediscriminante.R* et copier les lignes suivantes :

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../TP/TP/Code")
3 # Packages utilisés dans la suite
4 library(MASS)
5 # Supprimer toutes les variables
6 rm(list=ls(all=TRUE))
7 # Supprimer tous les graphiques déjà présents
8 graphics.off()
```

2. Lire les données.

```
1 # Lecture des données d'apprentissage
2 data_train <- read.table("../Data/synth_train.txt", header=T, sep="\t");
3 data_train_y <- data_train$y
4 data_train$y <- as.factor(data_train$y)
```

3. Appliquer la méthode d'analyse discriminante linéaire.

```
1 # Analyse discriminante LINEAIRE
2 lin_disc_an <- lda(y~., data = data_train)
```

4. Lire les données test et les prédire.

```
1 # Lecture des données test
2 data_test <- read.table("../Data/synth_test.txt", header=T, sep="\t");
3 data_test_x <- data.frame(x1=data_test$x1, x2=data_test$x2)
4 # Prédiction sur les données test
5 test_LDA_predict <- predict(lin_disc_an, newdata=data_test_x, type="class")
```

5. Évaluer la méthode sur le jeu de données test.

```
1 # Comparaison des valeurs prédites et des valeurs observées
2 table(test_LDA_predict$class, data_test$y)
3 # Calcul du taux d'erreur
4 lda_error_rate <- mean(test_LDA_predict$class != data_test$y)
5 cat("error rate using test data (LDA) = ", lda_error_rate)
```

6. Appliquer cette fois ci la méthode d'analyse discriminante quadratique.

```
1 # Analyse discriminante QUADRATIQUE
2 quad_disc_an <- qda(y~., data = data_train)
3 # Prédiction sur les données test
4 test_QDA_predict <- predict(quad_disc_an, newdata=data_test_x, type="class")
```

7. Évaluer la méthode sur le jeu de données test.

```
1 # Comparaison des valeurs prédites et des valeurs observées
2 table(test_QDA_predict$class, data_test$y)
3 # Calcul du taux d'erreur
4 qda_error_rate <- mean(test_QDA_predict$class != data_test$y)
5 cat("error rate using test data (QDA) = ", qda_error_rate)
```

8. Comparer les 2 méthodes en affichant leurs frontières sur une grille.

```
1 # Frontières LDA et QDA
2 gridx1 <- seq(from=min(data_train$x1), to=max(data_train$x1), length.out=40)
3 gridx2 <- seq(from=min(data_train$x2), to=max(data_train$x2), length.out=40)
4 grid <- expand.grid(x1 = gridx1, x2 = gridx2)
5 data_grid_x <- data.frame(x1=grid[,1],x2=grid[,2])
6
7 test_LDA_predict_grid <- predict(lin_disc_an, newdata=data_grid_x, type="class")
8 test_QDA_predict_grid <- predict(quad_disc_an, newdata=data_grid_x, type="class")
9
10 par(mfrow=c(1:2))
11 plot(data_train$x1,data_train$x2,col=data_train$y,pch=16)
12 par(new=T)
13 plot(data_grid_x$x1,data_grid_x$x2,col=test_LDA_predict_grid$class,pch=8,cex=0.5,
14      ann=FALSE)
15
16 plot(data_train$x1,data_train$x2,col=data_train$y,pch=16)
17 par(new=T)
18 plot(data_grid_x$x1,data_grid_x$x2,col=test_QDA_predict_grid$class,pch=8,cex=0.5,
19      ann=FALSE)
```

Bayésien naïf

1. Créer un fichier *classification-bayesiennaif.R* et copier les lignes suivantes.

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../TP/Code")
3 # Packages utilisés dans la suite
4 library(e1071)
5 # Supprimer toutes les variables
6 rm(list=ls(all=TRUE))
7 # Supprimer tous les graphiques déjà présents
8 graphics.off()
```

2. Lire les données d'apprentissage.

```
1 # Lecture des données d'apprentissage
2 data_train <- read.table("../Data/synth_train.txt",header=T,sep="\t");
3 data_train_y <- data_train$y
4 data_train$y <- as.factor(data_train$y)
```

3. Appliquer la méthode du Bayésien naïf.

```
1 # Bayésien Naïf
2 naiv_bayes <- naiveBayes(y~., data = data_train)
```

4. Lire les données test et prédire la sortie.

```
1 # Lecture des données test
2 data_test <- read.table("../Data/synth_test.txt",header=T,sep="\t");
3 data_test_x <- data.frame(x1=data_test$x1,x2=data_test$x2)
4 # Prédiction sur les données test
5 naiv_bayes_predict <- predict(naiv_bayes, newdata=data_test_x, type="class")
```

5. Évaluer la méthode sur le jeu de données test.

```
1 # Comparaison des valeurs prédites et des valeurs observées
2 table(naiv_bayes_predict, data_test$y)
3 # Calcul du taux d'erreur
4 naiv_bayes_error_rate <- mean(naiv_bayes_predict != data_test$y)
5 cat("error rate using test data (Naive bayes) = ", naiv_bayes_error_rate)
```

6. Afficher la frontière de cette méthode.

```
1 gridx1 <- seq(from=min(data_train$x1), to=max(data_train$x1), length.out=50)
2 gridx2 <- seq(from=min(data_train$x2), to=max(data_train$x2), length.out=50)
3 grid <- expand.grid(x1 = gridx1, x2 = gridx2)
4 data_grid_x <- data.frame(x1=grid[,1],x2=grid[,2])
5 naiv_bayes_predict_grid <- predict(naiv_bayes, newdata=data_grid_x, type="class")
6 plot(data_train$x1,data_train$x2,col=data_train_y,pch=16)
7 par(new=T)
8 plot(data_grid_x$x1,data_grid_x$x2,col=naiv_bayes_predict_grid,pch=8,cex=0.5,
9      ann=FALSE)
```

Régression logistique

1. Créer un fichier *classification-reglogistique.R* et copier les lignes suivantes.

```
1 # Adresse du dossier où vous travaillez
2 setwd("/Users/.../TP/Code")
3 # Supprimer toutes les variables
4 rm(list=ls(all=TRUE))
5 # Supprimer tous les graphiques déjà présents
6 graphics.off()
```

2. Lire les données d'apprentissage.

```
1 # Lecture des données d'apprentissage
2 data_train <- read.table("../Data/synth_train.txt",header=T,sep="\t");
3 data_train_y <- data_train$y
4 data_train$y <- as.factor(data_train$y)
```

3. Appliquer la méthode de régression logistique sur les données d'apprentissage.

```
1 # Régression logistique
2 glm_train <- glm(y~., data = data_train, family=binomial())
```

4. Lire les données test et prédire la sortie.

```
1 # Lecture des données test
2 data_test <- read.table("../Data/synth_test.txt",header=T,sep="\t");
3 data_test_x <- data.frame(x1=data_test$x1,x2=data_test$x2)
4 # Prédiction sur les données test
5 glm_train_predict <- predict(glm_train, newdata=data_test_x, type="response")
6 result_glm_train_predict <- (glm_train_predict > 0.5)+1
```

5. Évaluer la méthode sur le jeu de données test.

```
1 # Comparaison des valeurs prédites et des valeurs observées
2 table(result_glm_train_predict, data_test$y)
3 # Calcul du taux d'erreur
4 glm_error_rate <- mean(result_glm_train_predict != data_test$y)
5 cat("error rate using test data (Logistic regression) = ", glm_error_rate)
```

6. Afficher la frontière de cette méthode.

```
1 # Frontières LDA et QDA
2 gridx1 <- seq(from=min(data_train$x1), to=max(data_train$x1), length.out=50)
3 gridx2 <- seq(from=min(data_train$x2), to=max(data_train$x2), length.out=50)
4 grid <- expand.grid(x1 = gridx1, x2 = gridx2)
5 data_grid_x <- data.frame(x1=grid[,1],x2=grid[,2])
6 glm_train_predict_grid <- predict(glm_train,newdata=data_grid_x,type="response")
7 result_glm_train_predict_grid <- (glm_train_predict_grid > 0.5)+1
8 plot(data_train$x1,data_train$x2,col=data_train_y,pch=16)
9 par(new=T)
10 plot(data_grid_x$x1,data_grid_x$x2,col=result_glm_train_predict_grid,pch=8,
11      cex=0.5,ann=FALSE)
```