#### Analyse de données M1 Statistique et économétrie - 2012-2013 $V.\ Monbet$

#### Analyse Discriminante Décisionnelle

La collection des iris de Fisher est probablement l'une des plus célèbres dans le domaine de la reconnaissance de formes. Bien qu'ancienne, elle continue de faire référence dans le domaine. Le corpus consigne des mesures biométriques relevées sur des échantillons végétaux de type iris

- 1. longueur du sépale (en centimètres)
- 2. largeur du sépale (en centimètres)
- 3. longueur du pétal (en centimètres)
- 4. largeur du pétal (en centimètres)

Les 150 individus de la base sont répartis en trois classes équilibrées (50 individus dans chaque catégorie) correspondant à trois espèces florales : *iris setosa,iris versicolor* et *iris virginica*. La classe *setosa* est linéairement séparable des deux autres, alors que *versicolor* et *virginica* ne le sont pas entre elles.

Nous proposons dans ce TP de comparer plusieurs méthodes d'analyse discriminante décisionnelle; plus précisément, nous allons considérer différents estimateurs des lois des variables explicatives pour construire des règles de classement.

Remarque - Vous pouvez commencer par traiter toutes les question portant sur l'analyse univariée.

# 1 Analyse descriptive

1. Sous R, on peut charger les données par l'instruction suivante data(iris). Faites une analyse descriptive rapide des données. Discuter notamment le graphique suivant :

2. L'analyse de la variance (ANOVA) permet de tester l'effet de la variable discrête (Species) sur les variables continues. La variable continue qui correspond à la plus grande statistique de test de Ficher observée est la variable qui conduira à la meilleure régle de classement si on choisit de travailler en dimension 1. Aidez-vous des instructions suivantes pour déterminer les deux variables les plus discriminantes parmis les 4 variables continues.

```
var.names = c("Sepal.Length", "Sepal.Width", "Petal.Length", "Petal.Width")
nb_var = length(var.names)
for (k in 1:nb_var){
  print(var.names[k])
  print(anova(aov(iris[,k]~Species, data = iris)))
}
```

- 3. Densités marginales univariées
  - (a) Par la suite, nous allons chercher à estimer les lois (ou la loi jointe) de ces deux variables afin de construire une régle de classement. Il est donc intéressant de visualiser leurs histogrammes. Pour la première variable, on fait par exemple

- (b) Adapter ces commandes pour réaliser le même graphique pour la seconde variable la plus discriminante.
- (c) Superposer aux histogrammes les densités des lois de Gauss s'ajustant au mieux aux lois des variables dans chacune des classes. On peut utiliser la fonction dnorm pour calculer les densités. Commenter les graphiques.
- 4. On peut visualiser une estimation des lois jointes en utilisant la fonction bkde2D du package KernSmooth et la fonction contour.

5. Superposer les contours des lois de Gauss bivariées correspondantes. Commenter le graphique.

## 2 Régles de décisions sous une hypothèse de normalité

### 2.1 Modèle homoscédastique gaussien

#### 2.1.1 Cas univarié

- 1. Rappeler ce que veut dire le mot homoscédactique.
- 2. On va supposer que les espèces *virginica* et *versicolor* ont la même variance pour la longueur de pétale. Cette hypothèse vous semble t'elle raisonnable? Pourquoi (on peut par exemple utiliser la fonction var.test pour justifier la réponse)? Écrire (à la main!¹) la régle de décision associée à la variable longueur de pétale et qui permet de séparer les espèces *virginica* et *versicolor*.
- 3. Avant d'exécuter les commandes ci-dessous, pariez entre vous sur le résultat qui va s'afficher! Si vous comprenez ce que font ces lignes de code, vous devez tomber très près du résultat.

```
pv = NULL
for (k in 1:1000)
```

<sup>1. &</sup>quot;à la main" signifie ici en utilisant un papier et un crayon. En effet, il n'existe pas de routine R qui donne ce genre de réponse.

```
{pv[k] = var.test(sqrt(.26)*rnorm(50), sqrt(.26)*rnorm(50))$p.value}
sum(pv<.05)/1000
Interpréter le graphique suivant
vsim = NULL
for (k in 1:1000) {vsim[k] = var(sqrt(.26)*rnorm(50))}
hist(vsim)
abline(v=.26,col="red",lwd=1.5)</pre>
```

4. Tracer sur un même graphique les probabilités a posteriori pour ces deux espèces. Le graphique obtenu vous semble t'il cohérent? Pourquoi?

#### 2.1.2 Cas bivarié

1. On suppose maintenant que les espèces virginica et versicolor ont la même covariance pour le couple longueur de pétale et largeur de pétale. Estimer la covariance puis écrire (encore à la main) et développer (toujours à la main) la régle de décision associée au couple (longueur de pétale, largeur de pétale) et qui permet de séparer les espèces virginica et versicolor.

On peut vérifier les résultats obtenus à l'aide de la fonction 1da du package MASS.

- 2. Tracer sur un même graphique les contours des densités a posteriori pour ces deux espèces. Ajouter les nuages de points en choisissant des couleurs appropriées. Le graphique obtenu vous semble t'il cohérent? Pourquoi?
- 3. On veut classer de nouveaux iris n'appartenant pas à la base de données. Ajoutez les sur le graphique en utilisant la fonction points avec l'option pch=10. Utiliser la règle ci-dessus pour les classer. Calculer leurs probabilités a posteriori d'être dans chacune des deux classes.

	Long. pétale	Larg. pétale	Espèce
1	5.5	1.6	
2	6	2	
3	4	1.2	

## 3 Modèle hétéroscédastique gaussien

Un test de comparaison des variances montre que les variances de largeurs de pétales ne sont pas égales dans les classes *virginica* et *versicolor*. Vous pouvez le vérifier facilement à l'aide de la fonction var.test. On peut aussi réaliser un test basé sur la statistique de Wilks pour comparer les covariances :

```
species = iris$Species[51:150] ;
Y = cbind(iris[51:150,3],iris[51:150,4] )
ma = manova(Y ~ species)
summary(ma,test="W")
```

Reprendre les questions du cas gaussien homoscédastique bivarié sous l'hypothèse d'hétéroscédasticité. Classer de nouveau les trois iris du tableau ci-dessus.

On peut vérifier les résultats obtenus à l'aide de la fonction qda du package MASS.

```
fitq2 = qda(Species~Petal.Length + Petal.Width,data=iris)
    # analyse discriminante quadratique
```

## 4 Modèle non paramétrique - estimateurs à noyaux

- 1. Considérons le cas univarié et la variable longueur de pétales. Estimer les densités par classe par un estimateur à noyau en utilisant la fonction density. Tracer les densités a posteriori correspondantes. En déduire (à la main) une régle de classement.
- 2. Cas bivarié. Tracer les densités a posteriori.

library(MASS) # chargement des librairies

# 5 Comparaison de différentes méthodes dans le cadre multivarié

On utilise le package MASS pour l'analyse discriminante de Fisher (cas gaussien) et le package class pour l'analyse discriminante de Bayes naive (plus proches voisins).

1. Ajuster les modèles.

```
library(class) # pour kNN

fitq.disl=lda(Species~.,data=iris) # analyse discriminante linéaire
fitq.disq=qda(Species~.,data=iris) # analyse discriminante quadratique
```

fitq.knn=knn(iris[,1:4],iris[,1:4],iris\$Species,k=5) # k plus proches voisins

2. Estimer le risque de Bayes associé à chacune des règles en utilisant les instructions suivantes et en déduire quelle est la meilleure règle de classement au sens du risque de Bayes.

```
# erreur d'apprentissage
table(iris[,"Species"],predict(fitq.disl,iris)$class)
table(iris[,"Species"],predict(fitq.disq,iris)$class)
table(iris[,"Species"],knn(iris[,1:4],iris[,1:4],iris$Species,k=5))
```

3. Pouvez vous expliquer le résultat suivant :

```
table(iris[, "Species"], knn(iris[,1:4], iris[,1:4], iris$Species, k=1))
```

4. Dans la question précédente, le risque de Bayes est estimé par resubstitution. Cette méthode d'estimation de l'erreur conduit généralement à des résultats trop optimistes et il est préférable de travailler en validation croisée. En pratique on retire quelques individus du jeux de données typiquement 1/3), on ajuste les modèles à partir des individus restant puis on utilise les individus retirés pour tester les modèles, autrement dit pour estimer le risque de Bayes.

```
n = length(iris[,5])
n.test = 50
i.test = sample(1:n,n.test,replace=FALSE);
i.train = setdiff(1:n,i.test)

fitq.disl=lda(Species~.,data=iris,subset=i.train)
fitq.disq=qda(Species~.,data=iris,subset=i.train)
fitq.knn=knn(iris[i.train,1:4],iris[i.test,1:4],iris$Species[i.train],k=5)

table(iris[i.test,"Species"],predict(fitq.disl,iris)$class[i.test])
table(iris[i.test,"Species"],predict(fitq.disq,iris)$class[i.test])
table(iris[i.test,"Species"],fitq.knn)
```

On peut faire varier le nombre de plus proches voisins pour essayer d'améliorer les résultats. Quels sont les individus mal classés?

## 6 Et que peut-on faire avec SAS?

On peut bien sûr mener les mêmes analyses et réaliser les mêmes graphiques avec SAS.

## 6.1 Analyse discriminante univariée

- 1. Lire les données en utilisant le programme iris\_discriminant\_analysis.sas.
- 2. Tracer les histogrammes de la variable largeur de pétale conditionellement à chacune des espèces en utilisant les instructions suivantes :

```
pattern1 c=red
                        /*v=11
     pattern2 c=yellow /*v=empty*/;
     pattern3 c=blue
                        /*v=r1
     axis1 label=(angle=90);
     axis2 value=(height=.6);
     legend1 frame label=none;
     proc gchart data=iris;
        vbar PetalWidth / subgroup=Species midpoints=0 to 25
              raxis=axis1 maxis=axis2 legend=legend1 cframe=ligr;
     run;
3. Faire l'analyse discriminante dans le cas homoscédastique gaussien.
  title2 'Using Normal Density Estimates with Equal Variance';
  proc discrim data=iris method=normal pool=yes
                testdata=plotdata testout=plotp testoutd=plotd
                short noclassify crosslisterr;
     class Species;
     var PetalWidth;
  Puis tracer les densités et les densités a posteriori en utilisant les macros (voir
  graphiques_iris.sas).
  %plotden;
  %plotprob;
4. Dans le cas hétéroscédastique, on fait
  title2 'Using Normal Density Estimates with Unequal Variance';
  proc discrim data=iris method=normal pool=no
                testdata=plotdata testout=plotp testoutd=plotd
                short noclassify crosslisterr;
     class Species;
     var PetalWidth;
  run;
```

```
%plotden;
  %plotprob;
5. Si on utilise des estimateurs à noyau
  title2 'Using Kernel Density Estimates with Equal Bandwidth';
  proc discrim data=iris method=npar kernel=normal
                   r=.4 pool=yes
                testdata=plotdata testout=plotp
                   testoutd=plotd
                short noclassify crosslisterr;
     class Species;
     var PetalWidth;
  run;
  %plotden;
  %plotprob;
6. Et enfin avec des estimateurs des plus proches voisins,
  title2 'Using Kernel Density Estimates with Unequal Bandwidth';
  proc discrim data=iris method=npar kernel=normal
                   r=.4 pool=no
                testdata=plotdata testout=plotp
                   testoutd=plotd
                short noclassify crosslisterr;
     class Species;
     var PetalWidth;
  run;
  %plotden;
  %plotprob;
```