
TP 4 - DISCRIMINATION

UV : **SY09**

Branche : **Génie Informatique**

Filière : **Fouille de Données et Décisionnel**

Auteurs : **LU Han - HAMONNAIS Raphaël**

Table des matières

1	Programmation	2
1.1	Analyse discriminante	2
1.1.1	Apprentissage	2
1.1.2	Discrimination	3
1.2	Régression logistique	3
1.2.1	Apprentissage	3
1.2.2	Discrimination	5
1.2.3	Régression Logistique Quadratique	5
1.3	Vérification des fonctions	6
1.3.1	Analyse discriminante	6
1.3.2	Régression logistique	7
2	Application	8
2.1	Test sur les données simulées	9
2.1.1	Visualisation graphique des données	9
2.1.2	Visualisation « mathématique » des données via l'estimation des paramètres	10
2.1.3	Frontières de décisions	11
2.1.4	Calculs des taux d'erreur	11
2.2	Test sur données réelles	13
2.2.1	Données « Pima »	13
2.2.2	Données « Breast Cancer Wisconsin »	13
3	Challenge : données « Spam »	14
A	Fonctions pour l'analyse discriminante	15
A.1	Fonctions d'apprentissage <code>adq.app</code> , <code>adl.app</code> et <code>nba.app</code>	15
A.2	Fonction de discrimination <code>ad.val</code>	17
B	Fonctions pour la régression logistique	18
B.1	Fonction de calcul des probabilités à posteriori <code>postprob</code>	18
B.2	Fonction d'apprentissage <code>log.app</code>	18
B.3	Fonction de discrimination <code>log.val</code>	19
B.4	Fonction pour la régression logistique quadratique	19
C	Frontières de décision	20
C.1	Données <i>Synth1</i>	20
C.2	Données <i>Synth2</i>	21
C.3	Données <i>Synth3</i>	22

1. Programmation

1.1 Analyse discriminante

On se situe ici dans le cadre binaire : on considère $g = 2$ classes.

1.1.1 Apprentissage

A chaque modèle d'analyse discriminante correspond une fonction d'apprentissage :

- » `adq.app` pour l'analyse discriminante quadratique ;
- » `adl.app` pour l'analyse discriminante linéaire ;
- » `nba.app` pour le classifieur bayésien naïf.
- » **Remarque** : n'ayant pas vu que des fonctions « à compléter » (même si c'était clairement dit dans le sujet) étaient disponibles, nous avons implémenté nos propres fonctions dont vous trouverez le code en annexe A.

Prototype des fonctions `nba.app`, `adq.app` et `adl.app`

- » Objectif :
 - Apprendre les paramètres (proportions, vecteurs de moyennes et matrices de covariance des deux classes) du modèle considéré.
- » Paramètres en entrée :
 - `Xapp` : Les données d'apprentissage, de taille $n \times p$;
 - `zapp` : Les étiquettes des données d'apprentissage, de taille n ;
- » Sortie :
 - Une structure `params` contenant les paramètres du modèle :
 - `params$pi` : liste des proportions de chaque classe ;
 - `params$mu` : liste des vecteurs de moyennes de chaque classe ;
 - `params$sigma` : liste des matrices de covariance de chaque classe ;
- » Calcul des paramètres :
 - Proportions `pi` identiques pour les trois modèles, égales à $\frac{n_k}{n}$;
 - Vecteurs de moyennes `mu` identiques pour les trois modèles, égaux à $\bar{x}_k = \frac{1}{n_k} \sum_{i=1}^n z_{ik} x_i$ avec $z_{ik} = 1$ si $x_i \in \text{Classe } k$, 0 sinon ;
 - Matrices de covariance `sigma` différent selon les modèles :
 - AD Quadratique : chaque classe k a sa propre matrice de covariance V_k ;
 - AD Linéaire : la matrice de covariance est commune à toutes les classes et vaut $\Sigma = \frac{1}{n} \sum_{k=1}^g n_k V_k$;
 - Classifieur bayésien naïf : hypothèse d'indépendance conditionnelle des variables, on a alors Σ_k égale à la matrice diagonale obtenue à partir des valeurs diagonales de la matrice de covariance V_k , soit `sigma_k = diag(diag(V_k))` en code R ;

Le code des fonctions d'apprentissage est disponible en annexe A.1.

1.1.2 Discrimination

Prototype de la fonction de discrimination `ad.val`

- » Objectif :
 - Calculer les probabilités à postériori des individus de test pour chaque classe et effectuer le classement en fonction de ces probabilités.
- » Paramètres en entrée :
 - `params` : La structure contenant les paramètres du modèle utilisé pour effectuer la discrimination ;
 - `Xtst` : Les données de test à classer.
- » Sortie :
 - Une structure `discrimination` contenant les probabilités à postériori et le classement associé :
 - `discrimination$pw1` : les probabilités à postériori pour la classe 1 ;
 - `discrimination$pw2` : les probabilités à postériori pour la classe 2 ;
 - `discrimination$pred` : le classement des données `Xtst` dans la classe 1 ou 2 ;
- » Calcul des probabilités à postériori :
 - Utilisation de la formule $\mathbb{P}(\omega_k|x) = \frac{f_k(x)\pi_k}{f(x)} = \frac{f_k(x)\pi_k}{f_1(x)\pi_1 + f_2(x)\pi_2}$;
 - Utilisation de la fonction `mvdnorm` afin de calculer $f_1(x)$ et $f_2(x)$:
 - $f_1(x) = \text{mvdnorm}(Xtst, \mu_1, \Sigma_1)$
 - $f_2(x) = \text{mvdnorm}(Xtst, \mu_2, \Sigma_2)$
 - Remarquons que dans notre cas, avec 2 classes, $\mathbb{P}(\omega_2|x) = 1 - \mathbb{P}(\omega_1|x)$.
- » Détermination de la classe :
 - Si $P(\omega_1|x) \geq P(\omega_2|x)$, alors x appartient à la classe 1, sinon à la classe 2.

Le code de la fonction de discrimination est disponible en annexe A.2.

1.2 Régression logistique

1.2.1 Apprentissage

On rappelle qu'on est dans un cas binaire, c'est à dire avec $g = 2$ classes.

Alors

$$\mathbf{t} = (t_1, \dots, t_n)^t \quad \text{tel que} \quad t_i = \begin{cases} 1 & \text{si } Z_i = \omega_1 \\ 0 & \text{si } Z_i = \omega_2 \end{cases}$$

$$\mathbf{p} = p(x; \beta) = \mathbb{P}(\omega_1|x)$$

$$\mathbf{1-p} = 1 - p(x; \beta) = \mathbb{P}(\omega_2|x)$$

Notons pour plus de simplicité

$$\mathbf{p_i} = p(x_i; \beta) = \mathbb{P}(\omega_1|x_i)$$

$$\text{et } \mathbf{p} = (p_1, \dots, p_n)^t$$

Calcul des probabilités à postériori avec la fonction postprob

» Objectif :

- Calculer les probabilités à postériori de la classe 1 $\mathbb{P}(\omega_1|x)$.

» Paramètres en entrée :

- **beta** : estimateur du maximum de vraisemblance $\hat{\beta}$ des paramètres du modèle ;
- **X** : données dont on veut calculer $\mathbb{P}(\omega_1|x)$.

» Sortie :

- Le vecteur des probabilités $\mathbb{P}(\omega_1|x)$.

On utilise la formule suivante :

$$\mathbb{P}(\omega_1|x) = \frac{\exp(\beta^t x)}{1 + \exp(\beta^t x)}$$

Le code de la fonction **postprob** est disponible en annexe B.1.

Prototype de la fonction log.app

» Objectif :

- Apprendre les paramètres (matrice **beta**) du modèle de régression logistique à l'aide de l'algorithme de Newton-Raphson.

» Paramètres en entrée :

- **Xapp** : les données d'apprentissage, de taille $n \times p$;
- **zapp** : les étiquettes des données d'apprentissage, de taille n ;
- **intr** : variable binaire indiquant s'il faut ou non ajouter une ordonnée à l'origine (*intercept*) à la matrice **Xapp** ;
- **epsi** : scalaire correspondant au seuil ϵ en-deça duquel on considère que l'algorithme d'apprentissage a convergé.

» Sortie :

- Une structure **params** contenant les paramètres du modèle :
 - **params\$beta** : estimateur du maximum de vraisemblance $\hat{\beta}$ des paramètres, de dimensions $p \times 1$ (ou $(p+1) \times 1$ si une ordonnée à l'origine a été ajoutée) ;
 - **params\$niter** : le nombre d'itération effectuées par l'algorithme de Newton-Raphson avant convergence ;
 - **params\$logL** : valeur de la vraisemblance à l'optimum.

Calcul des paramètres β

Estimateur du maximum de vraisemblance $\hat{\beta}$ des paramètres :

» La méthode de Newton-Raphson consiste à approximer β par développement limité :

- On choisit un point de départ $\beta_{(q)}$, ($\beta_{(0)}$ par exemple) ;
- On calcule $\nabla \ln L_{(q)} = \frac{\partial \ln L(\beta_{(k)})}{\partial \ln L}$;
- On calcule $H_{(q)} = \frac{\partial^2 \ln L(\beta_{(k)})}{\partial \beta \partial \beta^t}$;
- On calcule $\beta_{(q+1)}$ en fonction de $\beta_{(q)}$, $\nabla \ln L_{(q)}$ et $H_{(q)}$;

- On itère jusqu'à la convergence définie par la valeur ϵ (c'est à dire $\nabla \ln L$ proche de zéro, donc maximum de vraisemblance atteint).
- » Équation de mise à jour de β
 - $\beta_{(q+1)} = \beta_{(q)} - H_{(q)} \nabla \ln L_{(q)}$;
 - Avec $\nabla \ln L_{(q)} = X^t(t - p_{(q)})$ où $p_{(q)} = p(x; \beta_{(q)})$
 - Avec $H_{(q)} = -X^t W_{(q)} X$ où $W_{(q)} = \text{diag}(p_{(q)}(1 - p_{(q)}))$

Valeur de la vraisemblance à l'optimum :

»

Le code de la fonction d'apprentissage `log.app` est disponible en annexe B.2.

1.2.2 Discrimination

Prototype de la fonction `log.app`

- » Objectif :
 - Classer les individus dans leur classes respectives en fonction des probabilités à postériori de chaque classe.
- » Paramètres en entrée :
 - `beta` : estimateur du maximum de vraisemblance $\hat{\beta}$ des paramètres du modèle ;
 - `Xtst` : données à discriminer de taille $n \times p$.
- » Sortie :
 - Une structure `discrimination` contenant les probabilités à postériori et le classement associé :
 - `discrimination$prob` : matrice de taille $n \times 2$ les probabilités à postériori des classes 1 et 2 ;
 - `discrimination$pred` : vecteur de taille n contenant le classement des données `Xtst` dans la classe 1 ou 2 ;
- » Détermination de la classe :
 - Si $P(\omega_1|x) \geq P(\omega_2|x)$, alors x appartient à la classe 1, sinon à la classe 2.

Le code de la fonction de discrimination `log.val` est disponible en annexe B.3.

1.2.3 Régression Logistique Quadratique

Cette méthode consiste à transformer l'espace des données en un espace plus complexe (contenant plus de variables) et à appliquer la régression logistique sur ce nouvel espace.

Le modèle ainsi défini est donc plus flexible (il permet de construire une frontière de décision plus complexe) ; mais le nombre de paramètres à estimer étant plus important, il peut également s'avérer moins robuste que le modèle classique déterminé dans l'espace des caractéristiques initiales.

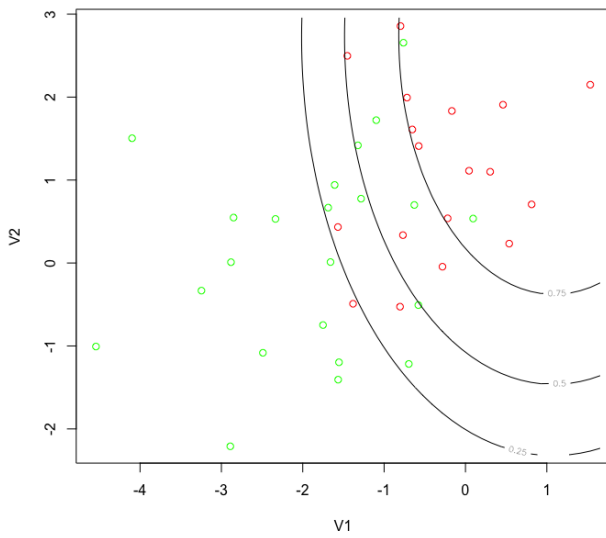
Par exemple, dans le cas où les individus sont décrits par les variables naturelles X_1, X_2 et X_3 , la régression logistique quadratique consiste à apprendre un modèle de régression logistique classique dans l'espace $\mathbf{X}^2 = \{X^1, X^2, X^3, X^1 X^2, X^1 X^3, X^2 X^3, (X^1)^2, (X^2)^2, (X^3)^2\}$, plutôt que dans l'espace $\mathbf{X} = \{X^1, X^2, X^3\}$.

La dimension du jeu de données p devient alors $p(p+3)/2$.

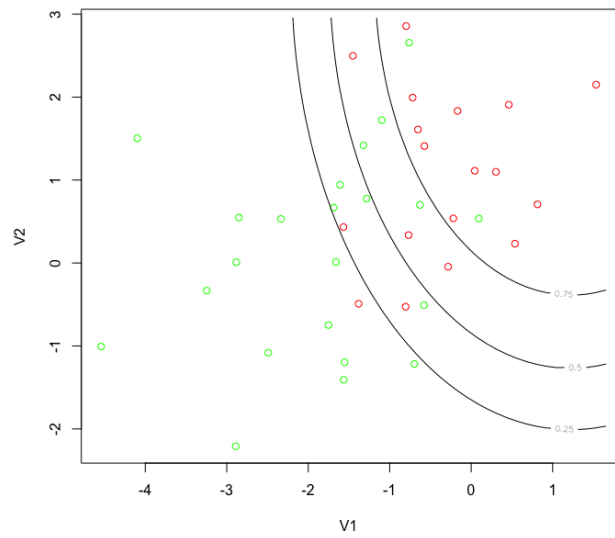
Le code de la fonction `log_quad` transformant l'espace originel en un espace quadratique est disponible en annexe B.4.

1.3 Vérification des fonctions

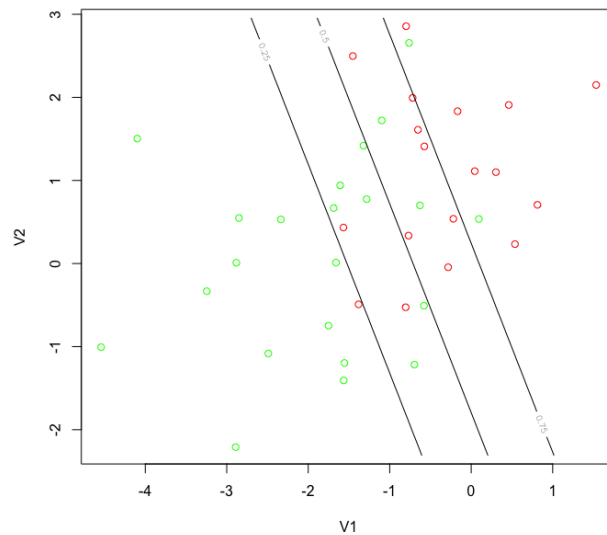
1.3.1 Analyse discriminante



(a) A. D. Quadratique



(b) Classifieur bayésien naïf



(c) A. D. Linéaire

FIGURE 1.1 – Frontières de décision pour l'analyse discriminante

1.3.2 Régression logistique

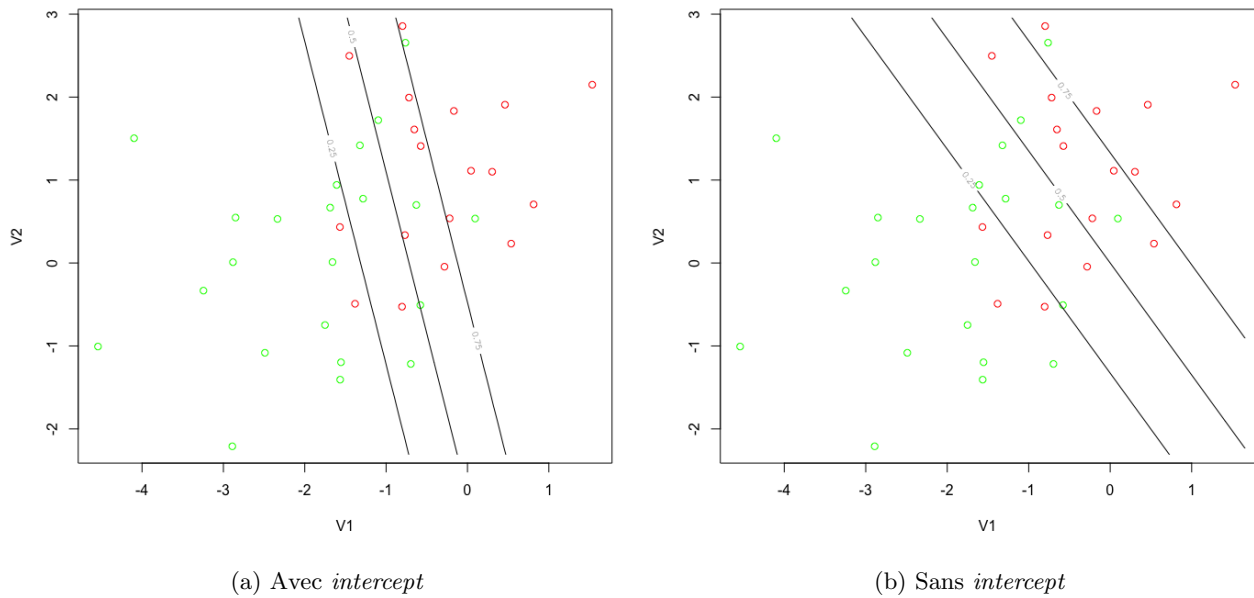


FIGURE 1.2 – Frontières de décision pour la régression logistique

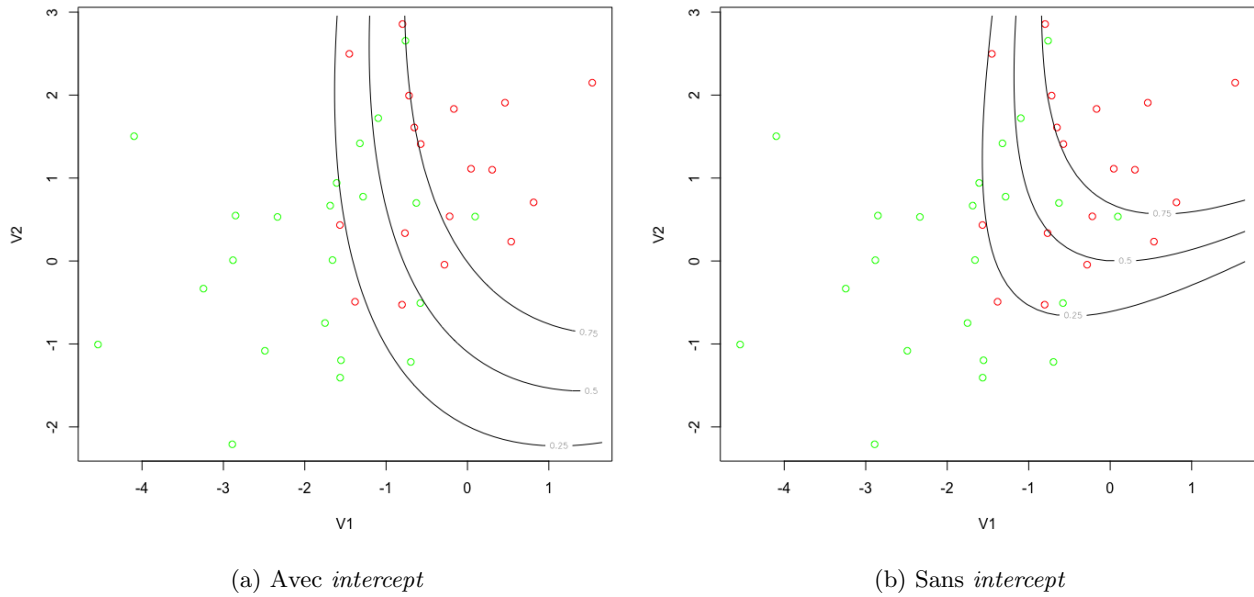


FIGURE 1.3 – Frontières de décision pour la régression logistique quadratique

2. Application

On va s'attacher ici à comparer les performances de nos cinq classifieurs :

- » Analyse discriminante quadratique, linéaire et bayésien naïf ;
- » Régression logistique simple et quadratique ;
- » Arbre de décision.

Il est important de bien noter les différence entre ces familles de classifieur, et plus particulièrement entre l'analyse discriminante et les deux autres.

L'analyse discriminante **suppose** que les données, conditionnellement à chaque classe ω_k , suivent une loi normale multidimensionnelle d'espérance μ_k et de variance Σ_k . En faisant ensuite différentes **hypothèses sur ces paramètres**, on obtient différentes expressions de la règle de Bayes, d'où l'on déduit différentes probabilités à postériori et donc différentes règles de décision en remplaçant les paramètres théoriques par leurs estimations. Dans le cas qui nous intéresse, seules les hypothèses sur les matrices de variances changent en fonction des classifieurs, et donc leurs estimations changent :

- » AD Quadratique : chaque classe k a sa propre matrice de covariance V_k ;
 - Ce classifieur sera plus fidèle aux les données d'apprentissage et donnera une frontière de décision quadratique ;
 - Il pourra donc être très efficace dans le cas de classes à distribution quadratique mais risque de faire du sur-apprentissage dans le cas où les données d'apprentissage ne sont pas représentatives des données réelles.
- » AD Linéaire : la matrice de covariance est commune à toutes les classes et vaut $\Sigma = \frac{1}{n} \sum_{k=1}^g n_k V_k$;
 - Ce classifieur fait une supposition forte qui peut l'amener à avoir de moins bonne performances que l'AD quadratique ;
 - Par contre, la frontière de décision obtenue étant linéaire, elle sera plus robuste et donc moins sensible au bruit des données (i.e. données d'apprentissage non représentatives).
- » Classifieur bayésien naïf : hypothèse d'indépendance conditionnelle des variables, on a alors Σ_k égale à la matrice diagonale obtenue à partir des valeurs diagonales de la matrice de covariance V_k .
 - Ce classifieur propose des frontières de décision quadratique, tout comme l'A.D. quadratique ;
 - Par contre, l'hypothèse d'indépendance conditionnelle des variables suppose que les classe sont des ellipsoïde dont les axes sont parallèles aux axes du repère.

Dans le cas où les données **respectent effectivement les hypothèses avancées** par les classifieurs d'analyse discriminante, alors il est fort probable que ces derniers présentent de bonnes performances. En effet, les estimations des probabilités à postériori seront d'autant plus précises que les hypothèses portant sur la distribution des données sont vérifiées.

Mais cela n'est pas obligatoire : des classes dont les espérances μ_k sont très proches seront dans la plupart des cas très difficiles à discriminer.

La régression logistique **ne fait pas d'hypothèses sur les distributions conditionnelles f_k** : elle estime directement les probabilités à postériori $\mathbb{P}(\omega_k|x)$ d'appartenance aux classes.

La régression logistique simple donnera une frontière de décision linéaire. Lorsqu'on change la dimension des données en rajoutant des variables qui sont des combinaisons des variables originales, la régression logistique devient quadratique et on pourra alors avoir un hyperplan comme frontière de décision.

2.1 Test sur les données simulées

On souhaite comparer les performances de l'analyse discriminante, de la régression logistique (linéaire et quadratique), et des arbres de décision sur les jeux de données simulées *Synth1*, *Synth2* et *Synth3*.

Pour ce faire, on répétera $N = 20$ fois pour chaque jeu de données :

- » séparer le jeu de données en un ensemble d'apprentissage et un ensemble de test ;
- » apprendre le modèle sur l'ensemble d'apprentissage ;
- » effectuer le classement des données de test et calculer le taux d'erreur associé.

Pour chaque jeu de données, on calculera ensuite le taux d'erreur (de test) moyen $\hat{\epsilon}$ sur les $N = 20$ séparations effectuées.

2.1.1 Visualisation graphique des données

Les données synthétiques sont de dimension $p = 2$, ce qui nous permet d'en tracer une visualisation graphique simple. Cette visualisation nous donnera des intuitions quand aux modèles de discrimination qui pourront être efficaces ou non sur ces données.

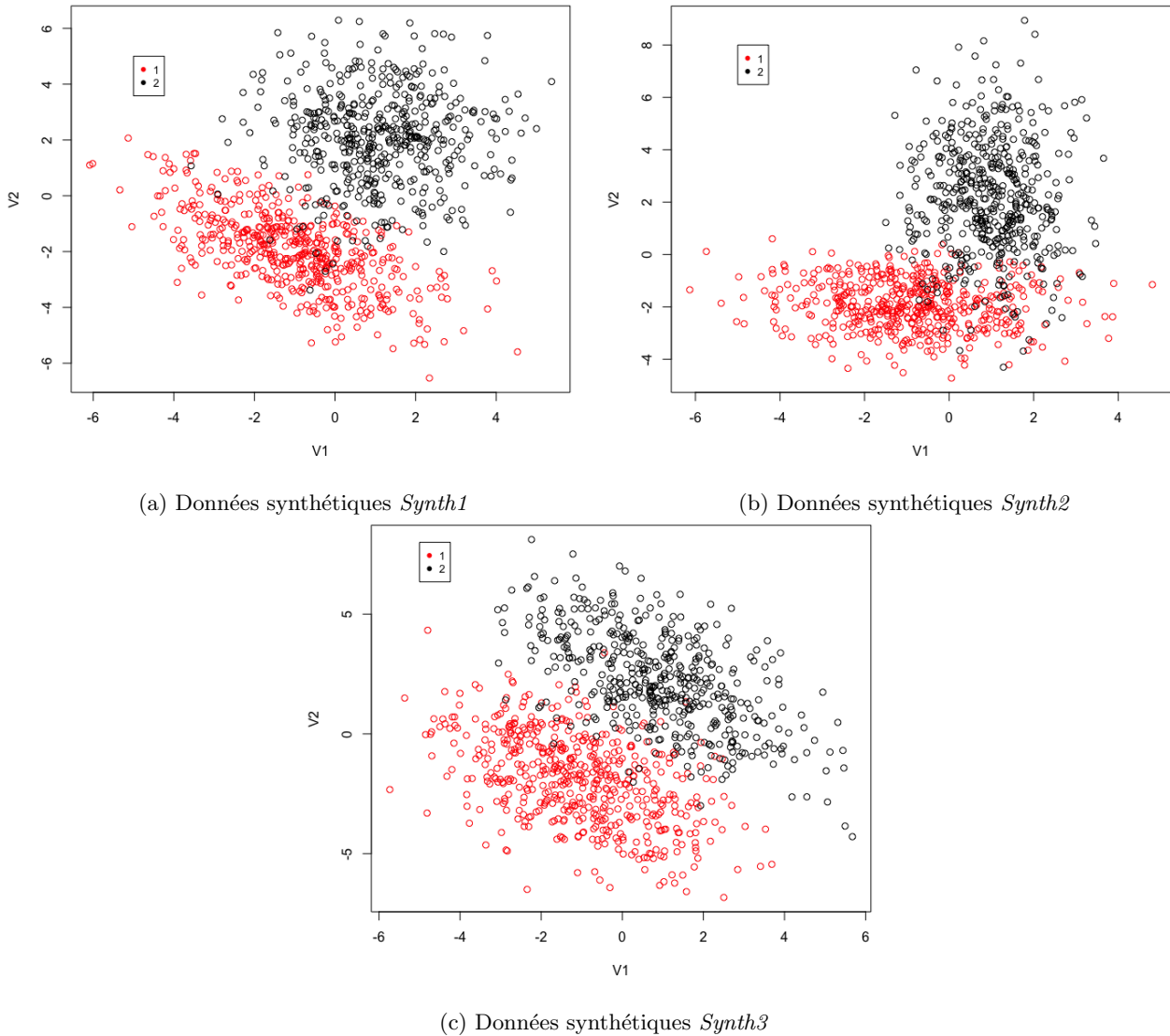


FIGURE 2.1 – Visualisation des données synthétiques *Synth1*, *Synth2* et *Synth3*

Remarques possibles suite à la visualisation :

» Données *Synth1* :

- Classes dont les variances Σ semblent différentes, Σ_1 donnant une forme d'ellipsoïde non parallèle aux axes quand Σ_2 semble définir une classe plutôt circulaire ;
- Il semble y avoir peu de chevauchement entre les classes ;
- Compte tenu de leurs formes, un classifieur à frontière de décision quadratique tendra à être plus efficace qu'un classifieur à frontière de décision linéaire.

» Données *Synth2* :

- Classes dont les variances Σ semblent inversées : la classe ω_1 est très étendue sur l'axe V1 quand la classe ω_2 est très étendue sur l'axe V2 ;
- Le chevauchement entre les classes est plus important que dans les deux autres jeux de données ;
- De façon plus prononcée que pour *Synth1* ou *Synth3*, un classifieur à frontière de décision quadratique sera plus performant qu'un classifieur à frontière de décision linéaire.

» Données *Synth3* :

- Classes dont les variances Σ semblent identiques ;
- Peu de chevauchement entre les classes ;
- L'analyse discriminante linéaire semble particulièrement indiquée ici. On peut même s'attendre à ce que les frontières de décision quadratiques aient une forme quasiment linéaire.

» Classifieur des arbres de décision :

- Les classes se présentant visuellement sous forme d'ellipsoïdes (pas toujours parallèles aux axes), ce classifieur – proposant une frontière de décision linéaire par morceaux **et** parallèle aux axes des dimensions (= variables) – pourra logiquement donner de plus mauvais résultats que les autres sur ces jeux de données.

2.1.2 Visualisation « mathématique » des données via l'estimation des paramètres

Fichier source	Centres de gravité μ	Variances Σ	Proportions π
Synth1	$\mu_1 = \begin{pmatrix} -1.03 \\ -1.95 \end{pmatrix}, \mu_2 = \begin{pmatrix} 1.10 \\ 2.07 \end{pmatrix}$	$\Sigma_1 = \begin{pmatrix} 2.89 & -1.53 \\ -1.53 & 1.96 \end{pmatrix}, \Sigma_2 = \begin{pmatrix} 2.09 & 0.00 \\ 0.00 & 2.81 \end{pmatrix}$	$\begin{pmatrix} \pi_1 = 0.51 \\ \pi_2 = 0.49 \end{pmatrix}$
Synth2	$\mu_1 = \begin{pmatrix} -0.90 \\ -1.92 \end{pmatrix}, \mu_2 = \begin{pmatrix} 0.97 \\ 2.17 \end{pmatrix}$	$\Sigma_1 = \begin{pmatrix} 2.81 & -0.19 \\ -0.19 & 0.90 \end{pmatrix}, \Sigma_2 = \begin{pmatrix} 0.90 & -0.04 \\ -0.04 & 4.60 \end{pmatrix}$	$\begin{pmatrix} \pi_1 = 0.5 \\ \pi_2 = 0.5 \end{pmatrix}$
Synth3	$\mu_1 = \begin{pmatrix} -1.06 \\ -1.91 \end{pmatrix}, \mu_2 = \begin{pmatrix} 0.92 \\ 2.19 \end{pmatrix}$	$\Sigma_1 = \begin{pmatrix} 2.88 & -1.55 \\ -1.55 & 3.50 \end{pmatrix}, \Sigma_2 = \begin{pmatrix} 2.92 & -2.02 \\ 2.02 & 4.17 \end{pmatrix}$	$\begin{pmatrix} \pi_1 = 0.52 \\ \pi_2 = 0.48 \end{pmatrix}$

TABLE 2.1 – Paramètres estimés des fichiers *Synth1*, *Synth2* et *Synth3*

Remarques :

» Données *Synth1* :

- Comme nous en avons eu l'intuition précédemment lors de la visualisation, on a bien $\Sigma_1 \neq \Sigma_2$ avec ω_1 légèrement plus dispersée sur l'axe V1 après rotation et ω_2 légèrement plus dispersée sur l'axe V2 (et non circulaire comme on l'avait supposé) ;
- Compte tenu des paramètres estimés, la classification discriminante quadratique sera plus indiquée que la linéaire ou le classifieur bayésien naïf, car leurs hypothèses ne sont pas vérifiées ;

- Par équivalence, la régression logistique quadratique sera plus indiquée que la régression logistique simple.
- » Données *Synht2* :
 - Compte tenu des paramètres estimés, le classifieur bayésien naïf sera plus indiqué que la discrimination linéaire ou quadratique, car il semble bien qu'on ait l'hypothèse d'indépendance conditionnelle des variables vérifiée (covariances nulles ou insignifiantes dans les matrices Σ_k);
 - Par équivalence, la régression logistique quadratique sera plus indiquée que la régression logistique simple (le classifieur bayésien naïf donne une frontière de décision quadratique).
- » Données *Synht3* :
 - Compte tenu des paramètres estimés, l'analyse discriminante linéaire semble la plus indiquée. Même si les matrices Σ_1 et Σ_2 ne sont pas exactement égales entre-elles, elles sont tout de même relativement proches;
 - Par équivalence, la régression logistique simple sera plus indiquée que la régression logistique quadratique.

2.1.3 Frontières de décisions

Des exemples de frontières de décisions sont disponibles en annexe :

- » Données *Synht1* : voir section C.1;
- » Données *Synht2* : voir section C.2;
- » Données *Synht3* : voir section C.3.

2.1.4 Calculs des taux d'erreur

Remarque sur l'espérance et la variance du taux d'erreur $\hat{\epsilon}$: on va par la suite calculer ces statistiques à des fins analytiques. Il faut prendre note que nous allons les calculer sur un échantillon de taille $n = 6$, voire moins lorsqu'on comparera les différents types de classifieur. Bien que donnant certaines indications et servant notre propos lors des comparaisons, ces statistiques doivent être considérées avec prudence.

Calcul du taux d'erreur $\hat{\epsilon}$ sur les données synthétiques :

Fichier source	ADQ	ADL	NBA	Log	Log Quad	Tree
Synth1	3.43	4.52	4.22	3.67	3.25	4.73
Synth2	6.23	8	6.17	6.85	6.44	8.09
Synth3	4.23	4.16	5.45	4.2	4.37	5.93

TABLE 2.2 – Estimation du taux d'erreur $\hat{\epsilon}$ sur $N = 20$ itérations des différents classifieurs sur les données synthétiques

Comme on avait pu le supposer précédemment lors de l'analyse des visualisations graphiques et des paramètres des classes, on a bien pour :

- » Données *Synht1* :
 - De meilleures performances avec l'analyse discriminante quadratique et la régression logistique quadratique.
- » Données *Synht2* :
 - De meilleures performances avec le classifieur bayésien naïf et la régression logistique quadratique.
- » Données *Synht3* :

— De meilleures performances avec l'analyse discriminante linéaire et la régression logistique simple.

» Classifieur des arbres de décision :

— Les plus mauvaises performance pour chaque jeu de données.

Il faut néanmoins nuancer nos propos car les performances des classifieurs par fichier sont à première vue sensiblement les mêmes (faible variance) :

Fichier source	Espérance de $\hat{\epsilon}$	Variance de $\hat{\epsilon}$
Synth1	3.97	0.37
Synth2	6.96	0.76
Synth3	4.72	0.59

TABLE 2.3 – Espérances et variances des taux d'erreurs de tous les classifieurs pour un fichier synthétique donné

Remarque : dans la suite de l'analyse des performances des classifieurs, on laissera de côté l'arbre de décision. Ce dernier propose une frontière linéaire parallèle aux axes « par morceaux » (voir frontières de décision en annexe Figure C.1f, Figure C.2f et Figure C.3f) et donne les résultats les plus mauvais, comparativement aux autres classifieurs.

Les données étant distribuées de façon ellipsoïdale, il n'est pas étonnant que ce soit ce classifieur qui présente la plus mauvaise performance.

Nous nous proposons maintenant de faire une analyse comparative des performances en deux parties :

- » les modèles d'analyse discriminante comparés à la régression logistique ;
- » les modèles à frontière de décision linéaire comparés aux modèles à frontière de décision quadratique.

Comparaison du taux d'erreur $\hat{\epsilon}$ entre analyse discriminante et régression logistique

Fichier source	AD : $\mathbb{E} [\hat{\epsilon}]$	AD : $Var [\hat{\epsilon}]$	RL : $\mathbb{E} [\hat{\epsilon}]$	RL : $Var [\hat{\epsilon}]$
Synth1	4.06	0.32	3.46	0.09
Synth2	6.8	1.08	6.64	0.08
Synth3	4.61	0.53	4.29	0.01

TABLE 2.4 – Espérances et variances des taux d'erreurs des classifieurs de type analyse discriminante versus les classifieurs de type régression logistique

Étonnamment, et quand bien même les données suivent dans chaque classe une loi normale multivariée et respectent à tour de rôle les hypothèses d'un classifieur d'analyse discriminante donné, on a quand même des performances *légèrement* supérieures pour la régression logistique (simple ou quadratique). Mais cette différence est tout de même trop faible pour être significative (0.35 d'écart en moyenne).

La variance du taux d'erreur est par contre significativement plus faible pour la régression logistique. Cela peut s'expliquer par le fait que cette dernière **ne fait pas de supposition** sur la distribution des données, contrairement à l'analyse discriminante.

Comparaison du taux d'erreur $\hat{\epsilon}$ entre classifieurs à frontière de décision linéaire ou quadratique

Fichier source	Front. Lin : $\mathbb{E} [\hat{\epsilon}]$	Front. Lin : $Var [\hat{\epsilon}]$	Front. Quad : $\mathbb{E} [\hat{\epsilon}]$	Front. Quad. : $Var [\hat{\epsilon}]$
Synth1	4.09	0.36	3.63	0.27
Synth2	7.42	0.66	6.28	0.02
Synth3	4.18	0	4.68	0.45

TABLE 2.5 – Espérances et variances des taux d'erreurs des classifieurs donnant une frontière de décision linéaire versus les classifieurs donnant une frontière de décision quadratique

Cette comparaison renforce les hypothèses émises précédemment lors de l'analyse des paramètres estimés des classes (voir sous-section 2.1.2) :

- » Les modèles quadratiques proposant des frontières de décision quadratiques sont plus performants sur les données *Synth1* et *Synth2* ;
- » Les modèles linéaires proposant des frontières de décision linéaires sont plus performants sur les données *Synth3*.

On remarquera aussi grâce à leur faibles variance que les estimateurs les plus adéquats sont aussi les plus fiables, c'est à dire les plus stables dans leurs performances.

2.2 Test sur données réelles

2.2.1 Données « Pima »

Fichier source	ADQ	ADL	NBA	Log	Log Quad	Tree
Pima	23.97	22.62	23.8	21.77	24.72	26.23

TABLE 2.6 – Estimation du taux d'erreur sur $N = 100$ itérations des différents classifieurs sur les données « Pima »

2.2.2 Données « Breast Cancer Wisconsin »

Fichier source	ADQ	ADL	NBA	Log	Tree
Breastcancer	4.93	4.46	3.99	3.94	5.39

TABLE 2.7 – Estimation du taux d'erreur sur $N = 100$ itérations des différents classifieurs sur les données « breast cancer Wisconsin »

3. Challenge : données « Spam »

A. Fonctions pour l'analyse discriminante

A.1 Fonctions d'apprentissage `adq.app`, `adl.app` et `nba.app`

Fonctions utilitaires de calcul des proportions, moyennes et matrices de co-variance

```
# calcul des proportions
prop.app = function(Xapp, zapp) {
  zapp = factor(zapp)
  g = length(levels(zapp)) # nombre de classes
  n = dim(Xapp)[1] # n = nombre d'individus total
  p = dim(Xapp)[2] # p = nombre de variables
  prop = list()
  for (k in 1:g) {
    prop[[k]] = table(zapp)[k] / length(zapp)
  }
  return(prop)
}

# calcul des centres de gravité
mu.app = function(Xapp, zapp) {
  zapp = factor(zapp)
  g = length(levels(zapp)) # nombre de classes
  n = dim(Xapp)[1] # n = nombre d'individus total
  p = dim(Xapp)[2] # p = nombre de variables
  mu = list()
  for (k in 1:g) {
    dataClassK = Xapp[zapp == levels(zapp)[k],]
    mu[[k]] = apply(dataClassK, MARGIN = 2, mean)
  }
  return(mu)
}

# calcul des matrices de covariance
sigma.app = function(Xapp, zapp) {
  zapp = factor(zapp)
  g = length(levels(zapp))
  n = dim(Xapp)[1]
  p = dim(Xapp)[2]
  sigma = list()
  for (k in 1:g) {
    dataClassK = Xapp[zapp == levels(zapp)[k],]
    sigma[[k]] = var(dataClassK)
  }
  return(sigma)
}
```


Fonctions d'apprentissage des paramètres

```

adq.app <- function(Xapp, zapp) {
  zapp = factor(zapp)
  params = list()
  params[["pi"]] = prop.app(Xapp, zapp)
  params[["mu"]] = mu.app(Xapp, zapp)
  params[["sigma"]] = sigma.app(Xapp, zapp)
  return(params)
}

adl.app <- function(Xapp, zapp) {
  zapp = factor(zapp)
  g = length(levels(zapp))
  p = dim(Xapp)[2]
  params = list()
  prop = prop.app(Xapp, zapp)
  mu = mu.app(Xapp, zapp)
  classes_sigma = sigma.app(Xapp, zapp)
  sigma = matrix(0, p, p)
  for (k in 1:g) {
    sigma = sigma + (prop[[k]] * classes_sigma[[k]])
  }
  for (k in 1:g)
    params[["sigma"]][[k]] = sigma
  params[["pi"]] = prop
  params[["mu"]] = mu
  return(params)
}

nba.app <- function(Xapp, zapp) {
  zapp = factor(zapp)
  g = length(levels(zapp))
  params = list()
  prop = prop.app(Xapp, zapp)
  mu = mu.app(Xapp, zapp)
  classes_sigma = sigma.app(Xapp, zapp)
  for (k in 1:g)
    classes_sigma[[k]] = diag(diag(classes_sigma[[k]]))
  params[["pi"]] = prop
  params[["mu"]] = mu
  params[["sigma"]] = classes_sigma
  return(params)
}

```

A.2 Fonction de discrimination ad.val

```
ad.val <- function(params, Xtst) {
  n = nrow(Xtst)

  f1 = mvdnorm(Xtst, params$mu[[1]], params$sigma[[1]])
  f2 = mvdnorm(Xtst, params$mu[[2]], params$sigma[[2]])

  discrimination = list()
  discrimination[["pw1"]] = vector(length = n)
  discrimination[["pw2"]] = vector(length = n)
  discrimination[["pred"]] = vector(length = n)

  for (i in 1 : n) {
    discrimination[["pw1"]][i] = f1[i]*params$pi[[1]]/(f1[i]*params$pi[[1]]+
f2[i]*params$pi[[2]])
    discrimination[["pw2"]][i] = f2[i]*params$pi[[2]]/(f1[i]*params$pi[[1]]+
f2[i]*params$pi[[2]])

    if(discrimination[["pw1"]][i] > discrimination[["pw2"]][i])
      discrimination[["pred"]][i] = 1
    else
      discrimination[["pred"]][i] = 2
  }
  discrimination[["pred"]] = factor(discrimination[["pred"]])
  return(discrimination)
}
```

B. Fonctions pour la régression logistique

B.1 Fonction de calcul des probabilités à postériori postprob

```
postprob <- function(beta, X) {  
  X <- as.matrix(X)  
  
  prob <- t(exp(t(beta)%*%t(X))/(1+exp(t(beta)%*%t(X))))  
}
```

B.2 Fonction d'apprentissage log.app

```
log.app <- function(Xapp, zapp, intr, epsi) {  
  n <- dim(Xapp)[1]  
  p <- dim(Xapp)[2]  
  
  Xapp <- as.matrix(Xapp)  
  
  if (intr == T) { # on ajoute une ordonnée à l'origine  
    Xapp <- cbind(rep(1,n),Xapp)  
    p <- p + 1  
  }  
  
  targ <- matrix(as.numeric(zapp),nrow=n) # ti: la réalisation d'une variable  $T_i \sim B(\pi_i)$   
  targ[which(targ==2),] <- 0 # remplacer la classe 2 par 0  
  Xapp_transposed <- t(Xapp)  
  
  beta <- matrix(0,nrow=p,ncol=1)  
  
  conv <- F  
  iter <- 0  
  while (conv == F) {  
    iter <- iter + 1  
    beta_old <- beta  
  
    prob_w1 <- postprob(beta, Xapp) #  $P(w_1|x)$   
    prob_w2 <- 1 - prob_w1  
    MatW <- diag(as.numeric(prob_w1 * prob_w2)) #  $W: W_{ii} = \pi_i(1-\pi_i)$   
  
    mat_hessienne <- -Xapp_transposed %*% MatW %*% Xapp  
    mat_hessienne_inverse <- solve(mat_hessienne)  
    gradient_w1 <- Xapp_transposed %*% (targ - prob_w1)  
    beta <- beta_old - (mat_hessienne_inverse %*% gradient_w1)  
  
    if (norm(beta - beta_old) < epsi) {  
      conv <- T  
    }  
  }  
  
  prob_w1 <- postprob(beta, Xapp) #  $P(w_1|x)$   
  prob_w2 <- 1 - prob_w1  
  out <- NULL  
  out$beta <- beta  
  out$iter <- iter  
  out$logL <- sum(targ*prob_w1+(1-targ)*(prob_w2))  
  out  
}
```

B.3 Fonction de discrimination log.val

```
log.val <- function(beta, Xtst) {
  m <- dim(Xtst)[1]
  p <- dim(beta)[1]
  pX <- dim(Xtst)[2]

  Xtst <- as.matrix(Xtst)

  if (pX == (p-1))
  {
    Xtst <- cbind(rep(1,m),Xtst)
  }

  prob_w1 <- postprob(beta, Xtst) # P(w1|x)
  prob_w2 <- 1 - prob_w1
  prob <- cbind(prob_w1, prob_w2)
  pred <- max.col(prob)

  out <- NULL
  out$prob <- prob
  out$pred <- pred
  return(out)
}
```

B.4 Fonction pour la régression logistique quadratique

Fonction qui transforme l'espace de X en un espace quadratique :

```
log_quad <- function(X) {
  X2 <- X

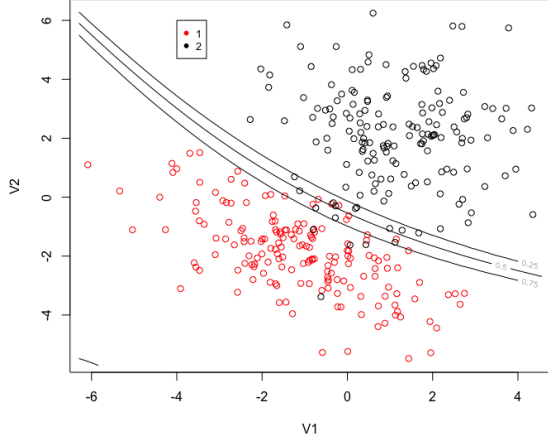
  for (p in 1:(dim(X)[2]-1))
    for (q in (p+1):dim(X)[2])
      X2 <- cbind(X2, X[,p]*X[,q])

  for (p in 1:dim(X)[2])
    X2 <- cbind(X2, X[,p]^2)

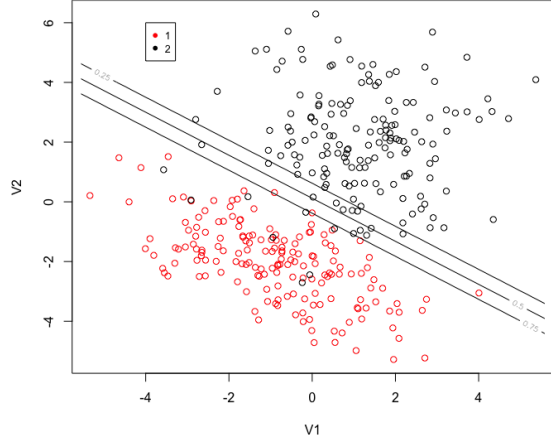
  return(X2)
}
```

C. Frontières de décision

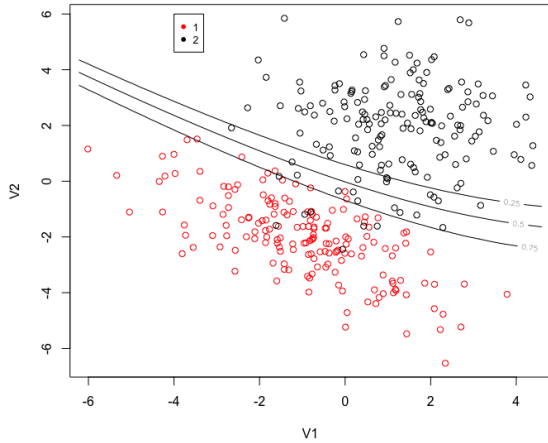
C.1 Données *Synth1*



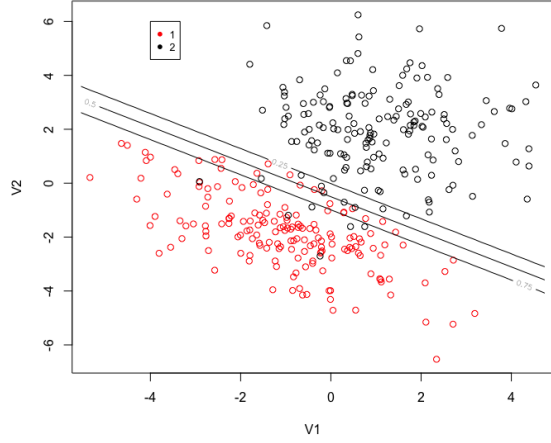
(a) A.D. Quadratique



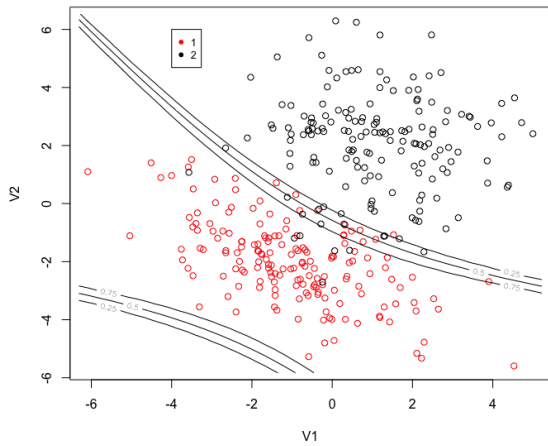
(b) A.D. Linaire



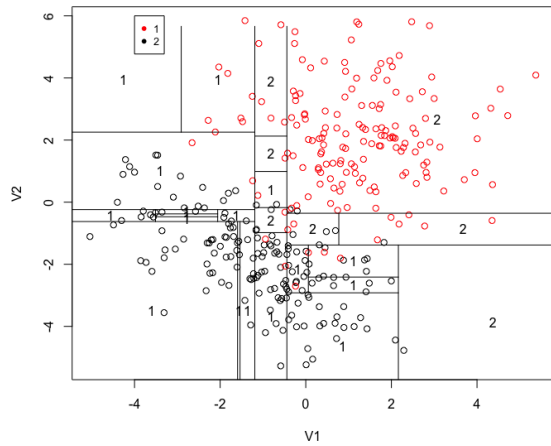
(c) Bayésien naïf



(d) Régression logistique



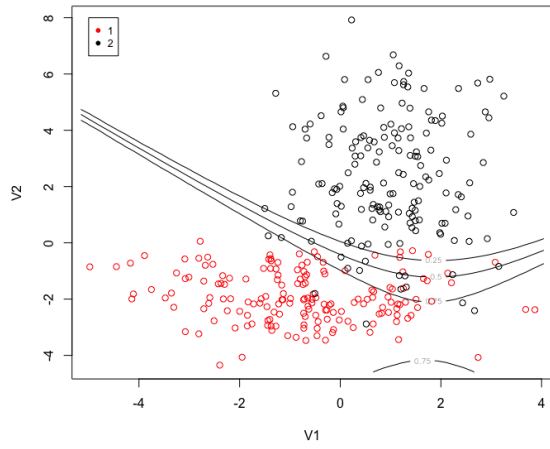
(e) Régression logistique quadratique



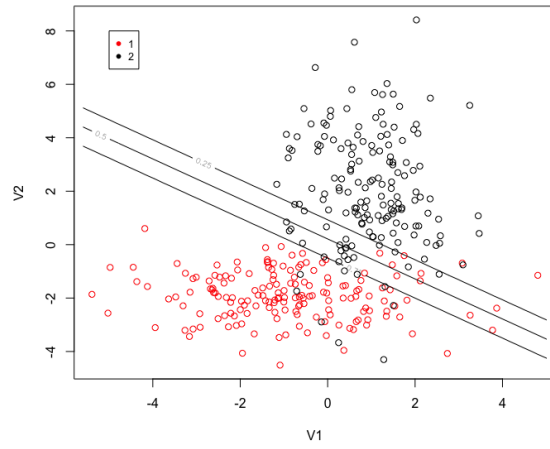
(f) Arbre de décision

FIGURE C.1 – Frontières de décision des classifieurs pour les données *Synth1* au niveaux $\mathbb{P}(\omega_1|x) = \{0.25, 0.5, 0.75\}$

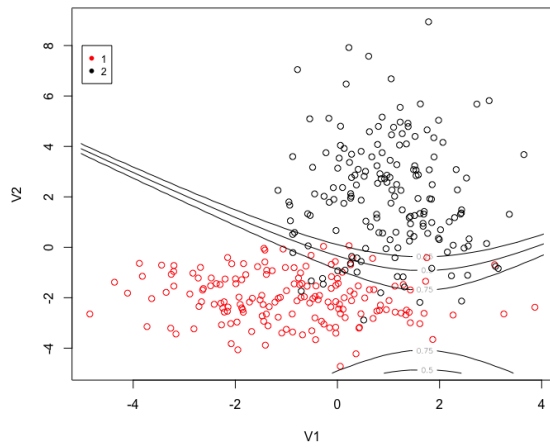
C.2 Données *Synth2*



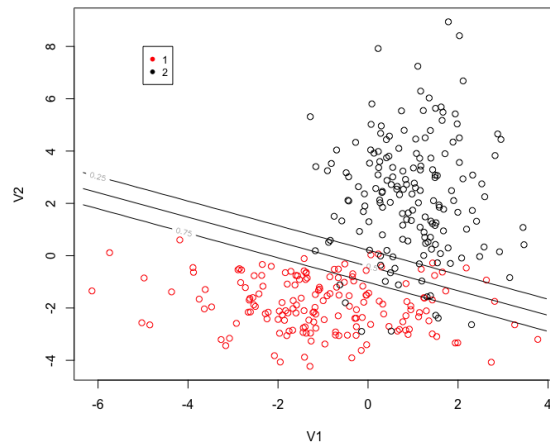
(a) A.D. Quadratique



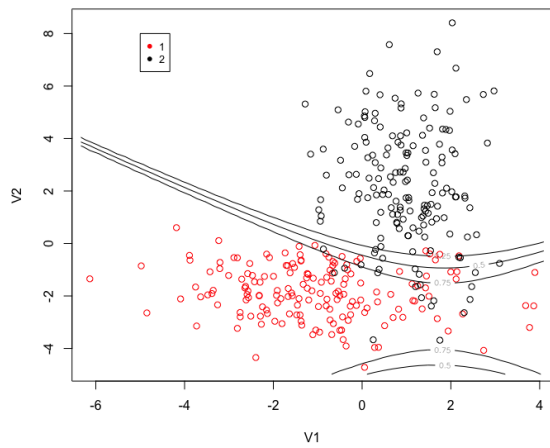
(b) A.D. Linaire



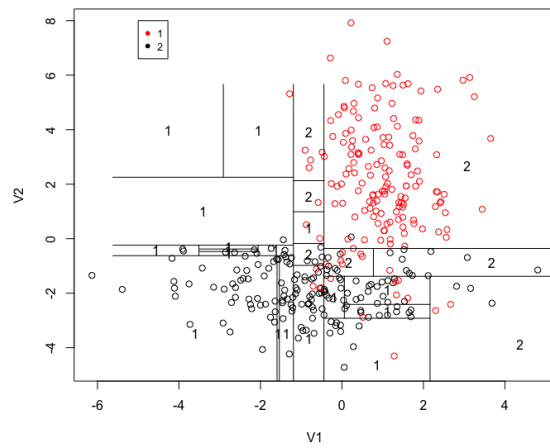
(c) Bayésien naïf



(d) Régression logistique



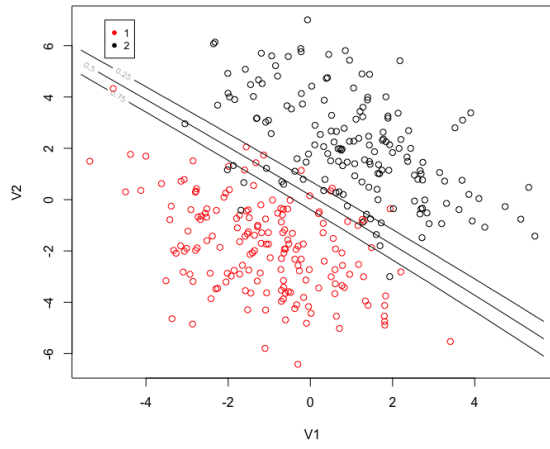
(e) Régression logistique quadratique



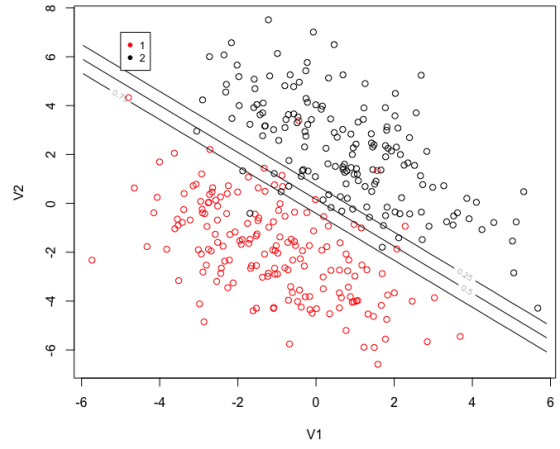
(f) Arbre de décision

FIGURE C.2 – Frontières de décision des classifieurs pour les données *Synth2* au niveaux $\mathbb{P}(\omega_1|x) = \{0.25, 0.5, 0.75\}$

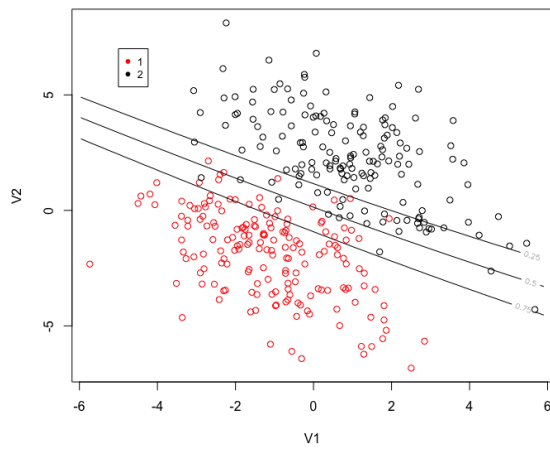
C.3 Données *Synth3*



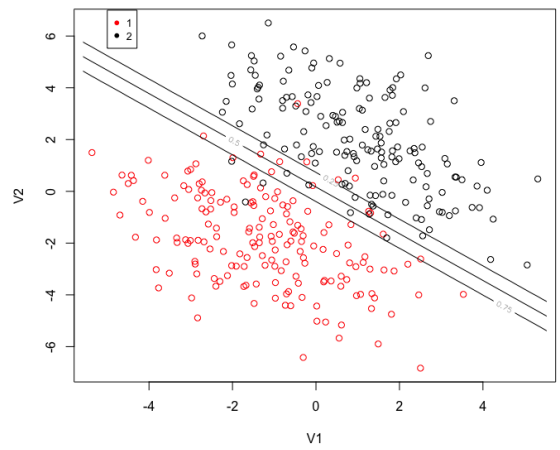
(a) A.D. Quadratique



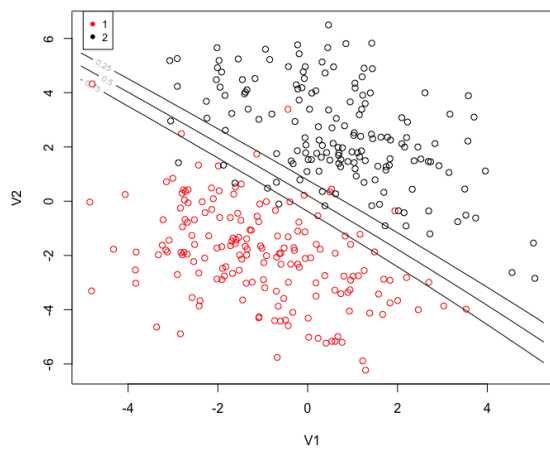
(b) A.D. Linaire



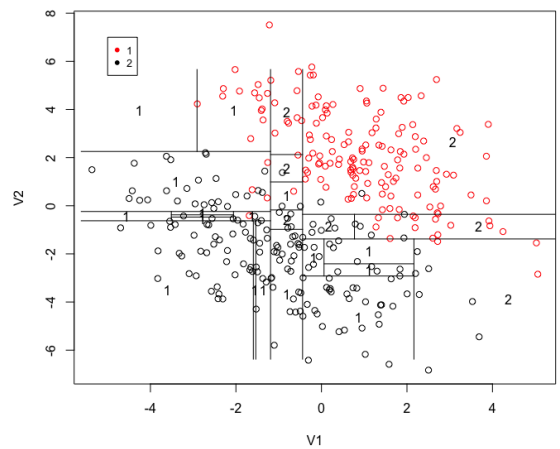
(c) Bayésien naïf



(d) Régression logistique



(e) Régression logistique quadratique



(f) Arbre de décision

FIGURE C.3 – Frontières de décision des classifieurs pour les données *Synth3* au niveaux $\mathbb{P}(\omega_1|x) = \{0.25, 0.5, 0.75\}$