

## 1 Données

Dans ces travaux dirigés on va analyser des données de puces avec la classification et la discrimination. Télécharger les données:

```
> library(golubEsets)
```

Il y a trois ensembles de données: Golub\_Train, Golub\_Test, and Golub\_Merge dans une class spéciale.

```
> data(Golub_Merge)
> ls()
```

Regardez les détails des données, on utilisera Golub\_Train comme ensemble d'apprentissage, et Golub\_Test comme ensemble test. Golub\_Merge pour faire l'analyse en validation croisée. et Golub\_Merge pour la classification.

## 2 Normalization of data

On regarde la distribution des données. On peut observer les données aberrantes. Il y a des nombres négatifs aussi à cause d'un enlèvement de fonds trop fort.

```
> hist(exprs(Golub_Merge))
```

Plus de détails:

```
> hist(exprs(Golub_Merge),breaks=100)
```

ou

```
> plot(density(exprs(Golub_Merge)))
```

Information sur l'expérience:

```
> sampleNames(Golub_Merge)
> phenoMerge <- pData(phenoData(Golub_Merge))
```

Faites le même pour les deux autres données (Golub\_Train and Golub\_Test) et comparez les résultats. Nous procédons aux données comme ont fait les auteurs de l'article, on modifie les valeurs extrêmes, les valeurs plus petites que 100 sont fixées à 100 et les valeurs plus grandes que 16,000 fixées à 16,000.

```
> gme <- exprs(Golub_Merge)
> gme[gme<100] <- 100
```

Faites pareils pour les valeurs grandes.

Maintenant prenons le log base 2 des données.

```
gme <- log2(gme)
```

Regardons les distributions de nouveau. Faites la meme chose avec les deux autres ensembles: on les appellera les matrices `gtre` et `gtee` a partir des données d'apprentissage et test.

Maintenant les données sont en bonne forme, on peut faire l'étape de filtrage.

### 3 Filtrage des genes

Une des des premieres taches consiste a enlever les genes qui ne seront pas interessants.

Moins de 40% des genes seront d'interet. On commence par enlever les genes qui ne sont pas exprimes et aussi les genes de maintenance qui sont exprimes de facon constant (housekeeping).

On utilise la bibliotheque `genefilter`.

```
>library(genefilter)
```

Les fonctions de filtrage lisent chaque ligne et produit une sortie booleenne qu'on peut utiliser pour prendre des sous ensembles de genes.

Essayer ce code:

```
mmfilt <- fonction(r = 5, d = 500, na.rm = TRUE){
  fonction(x) {
    minval <- min(2^x, na.rm = na.rm)
    maxval <- max(2^x, na.rm = na.rm)
    (maxval/minval > r) && (maxval - minval > d)
  }
}
```

Que fait cette fonction?

```
>ffun <- mmfilt()
>sub<- genefilter(gtre, ffun)
>gtres <- gtre[sub, ]
>gtees <- gtee[sub, ]
>gmes <- gme[sub, ]
```

Combien de genes restent -il?

### 3.1 Filtrage des genes

Il y a encore trop de genes qui restent?

Comme on veut etudier la difference de deux types de Lymphome on va se concentrer a cce sous ensemble. (ALL, AML).

On va garder les genes qui presentent une difference entre les deux groupes.

```
ttfun <- function(x){
  the.test <- t.test(x~phenoMerge$ALL.AML)
  if (the.test$p.value<10^(-4)) answer <- T
  else answer <- F
  return(answer)
}
```

La fonction ttfun utilise la fonction interne t.test pour faire le filtrage. On utilise le niveau  $p = 1e - 04$ . Nous utilisons l'information dans l'objet phenoMerge sur les types de leucemie.

Repetez ce que vous avez fait en utilisant la fonction ttfun (modifié pour utiliser les données phenoTrain et phenoTest) et les donnees gmes, gtres, et gtees. Appeler les resultats gmess, gtrss, et gteess.

Combien de genes restent-ils?

Pourquoi n'utilise-t-on pas  $p = 0.05$  comme la valeur de signification?

## 4 Location et Echelle pour standardization:

Standardization de location et echelle, on va centrer et standardiser les données: moyennes 0 et var 1.

```
star <- function(x) (x - mean(x))/sd(x)
merged <- t(apply(gmess, 1, star))
training <- t(apply(gtrss, 1, star))
testing <- t(apply(gteess, 1, star))
```

## 5 Discrimination

On utilise les methodes supervisees pour les categories ALL.AML. On utilise les données d'apprentissage pour construire le modele, on peut utiliser les methodes des voisins les plus proches pour assigner la classe. library: class.

On peut aussi utiliser les voisins les plus proches *knn()* pour classifier les patients Utiliser knn() avec les parametres  $k = 1, 11, 22, 33$  pour faire des classifications separees.

Le parametre cl: est la verification veritable des données d'apprentissage: phenoMerge\$ALL.AML.

Remarque: Nos matrices ont des patients comme colonnes et les variables(genes) comme lignes.

On doit transposer les données avant d'utiliser knn(). La fonction de transposition est t().

Quel est le taux de mauvais classement des différents essais? Quel est le meilleur choix de  $k$ ?

```
knn(t(training),t(testing),phenoTrain$ALL.AML) == phenoTest$ALL.AML
```

On utilise les valeurs d'expression dans les données assemblées et la fonction knn.cv .

```
knn.cv(t(merged),phenoMerged$A,k=1)==phenoMerged$A
```

Quel est le taux d'erreur?

## 6 Cluster Analysis:

On veut faire un groupement de patients on peut utiliser les distances euclidiennes calculées sur la matrice d'ensemble:

```
> merged.dist <- dist(t(merged))
```

(We need to transpose the data matrix, because dist() computes distances between the rows of the matrix. What happens if you skip the t() ?

Le résultat est une matrice 72 by 72. D'autres types de distances se calculent avec d'autres arguments, utiliser help(dist).

Créer une classification hiérarchique hierarchical clustering tree ou dendrogram.

```
> merged.ht <- hclust(merged.dist)
```

Faites help(hclust) pour voir les options.

L'objet de sortie de hclust s'utilise de plusieurs façons. La sortie graphique standard s'obtient avec:

```
> plot(merged.ht)
```

(plclust donne quelque chose de semblable).

Où coupez-vous?

Couper chaque arbre en deux groupes en utilisant la fonction cutree. Utiliser les deux groupes pour classer les patients:

```
> merged.hc.cut <- cutree(merged.dist,k=2)
> merged.hc.cut
> cbind(merged.hc.cut,phenoMerged$A)
```

Quel est le taux d'erreur?

Faites le graphique de classification avec la fonction heatmap.

## 7 Fonction de discrimination linéaire de Fisher

La méthode de discrimination linéaire (LDA) utilise les données d'apprentissage pour choisir les variables qui produisent une bonne "discrimination" des groupes.

Ensuite on peut essayer sur les nouvelles données.

```
library(MASS)
my.lda <- lda(t(training),groups=phenoTraining$A)
names(my.lda)
```

Maintenant on utilise la fonction de prediction sur les données tests:

```
my.predictions <- predict(my.lda,t(testing))
```

Quel est le taux d'erreur?

```
mean(my.predictions != phenoTesting$A)
```

Comparer à l'erreur de prediction estime par la méthode de validation croisée:

```
ldawithcv <- lda(t(training),groups=phenoTraining$A,CV=TRUE)
```