

Rapport Final statistique en grande dimension

AUTHOR

Amadou BAH ~ Frederic AKADJE ~ MARIA KONE

1. Chargement et exploration du jeu de données

Chargement du jeu de données

```
# Exemple : jeu de données simulé
set.seed(12311)
n <- 72
p <- 1000
X <- matrix(rnorm(n * p), n, p)
colnames(X) <- paste0("V", 1:p)
y <- factor(sample(c("ALL", "AML"), n, replace = TRUE))
#
# Vérifier les dimensions
dim(X)
```

```
[1] 72 1000
```

Exploration des données

Résumé des dimensions du jeu de données :

Compte tenu du nombre important de variables ici, une disposition particulière du jeu de données sera adoptée pour la réalisation de la description du jeu de données.

```
library(tidyverse)
library(gt)
# Transformation du jeu de données en dataframe
df <- bind_cols(X, Y=y)
# Nombre de valeurs manquantes par variable
na_ <- df |>
  select(-Y) |> # Retrait de la colonne cible
  summarise_all(list(na = ~sum(is.na(.)))) |> # Agrégation sur chaque colonne
  t() # Transposition (format vecteur)

# Moyenne des variables
mean_ <- df |>
  select(-Y) |>
  summarise_all(list(mean = ~mean(.x, na.rm = TRUE))) |>
  t()

# Maximum par colonne
max_ <- df |>
  select(-Y) |>
  summarise_all(list(max = ~max(.x))) |>
  t()
```

```

# Médiabe par colonne
median_ <- df |>
  select(-Y)|>
  summarise_all(list(max = ~max(.x))) |>
  t()

# Minimum par colonne
min_ <- df |>
  select(-Y)|>
  summarise_all(list(min = ~min(.x))) |>
  t()

# Ecart-type par colonne
sd_ <- df |>
  select(-Y)|>
  summarise_all(list(sd = ~sd(.x))) |>
  t()

# Tableau de restitution pour la description
restitution <- bind_cols(Variables = paste0("V", 1:p), # Concaténation de toutes les agrégatio
  NA_ =na_ ,
  Max_ = max_,
  Min_ = min_,
  Mean_ = mean_,
  Median_ = median_,
  SD_ = sd_ )

restitution |>
  head() |> # Affichage des 5 premières lignes
  gt() # Pour un format plus esthétiques

```

Variables	NA_	Max_	Min_	Mean_	Median_	SD_
V1	0	2.023186	-2.411731	-0.005736058	2.023186	1.0196627
V2	0	1.596161	-2.501312	-0.092842688	1.596161	0.9183733
V3	0	2.680843	-2.583187	-0.219670794	2.680843	1.1373863
V4	0	3.105307	-2.711703	0.046005704	3.105307	1.1202481
V5	0	2.693149	-1.999835	-0.137228616	2.693149	0.9328716
V6	0	1.803840	-1.675661	0.057810442	1.803840	0.8987529

Répartition des classes :

```

df |>
  group_by(Y) |>
  summarise(N=n(),`%`=round(n()/N*100,n,1)) |>

```

```
arrange(desc(N)) |>
gt()
```

	Y	N	%
AML	42	58.3	
ALL	30	41.7	

Deux classes sont représentées dans le jeu de données à savoir : ALL et AML.

Normalisation et échelle des variables ? :

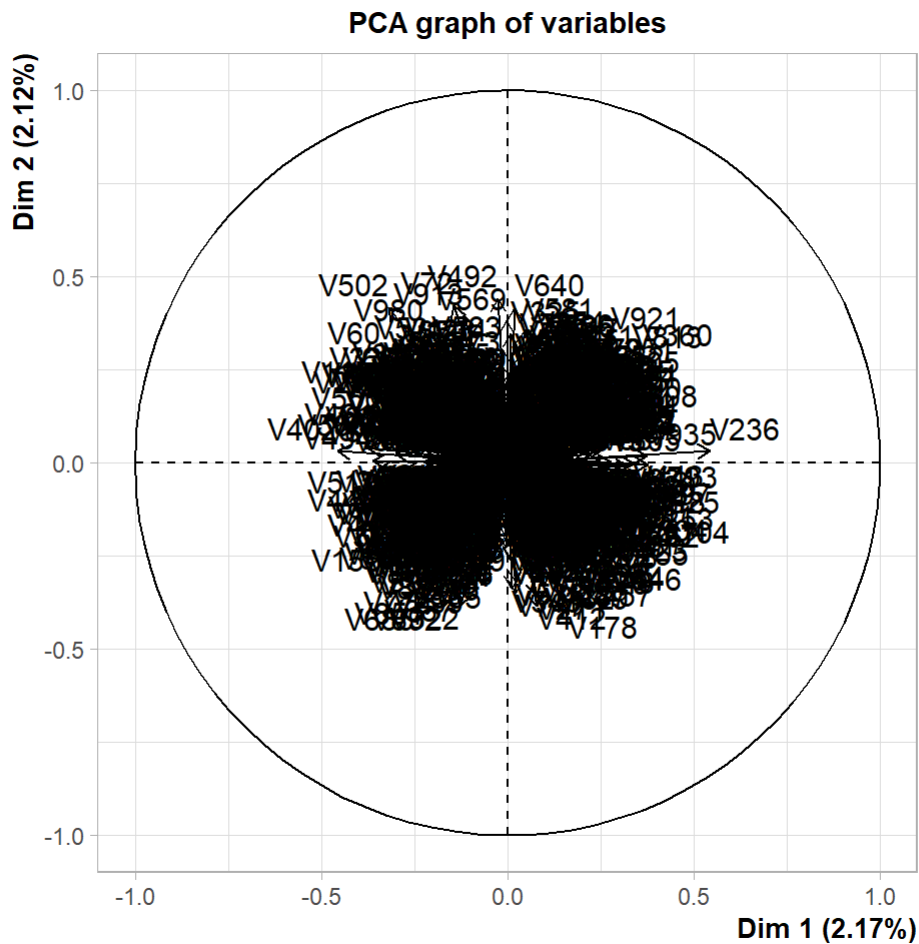
```
restitution |> summarise(
  `Valeurs manquantes`=sum(NA_),
  `Maximum des cols`=max(Max_),
  `Minimum des cols`= min(Min_),
  `Moyenne des cols` = mean(SD_),
  `Etendue des cols` = max(Max_)-min(Min_)
) |> gt()
```

Valeurs manquantes	Maximum des cols	Minimum des cols	Moyenne des cols	Etendue des cols
0	4.16039	-4.349511	0.9954469	8.509902

Le jeu de données comporte aucune valeur manquante au sein de chaque colonne. Toutefois, l'échelle des variables diffère à l'observation de l'étendue calculée à partir de l'ensemble des variables.

2. Analyse en Composantes Principales (ACP)

```
library(FactoMineR)
# ACP normée sur les 100 premières composantes principale
pcaRes <- PCA(df,scale.unit = TRUE,ncp = 100,quali.sup = 1001, graph = FALSE)
# Visualisation du premier plan principal
plot.PCA(pcaRes,choix = c('var'),axes = c(1,2))
```



A l'observation du graphique, les variables restent difficilement séparables. 4 grands regroupement se distinguent sur le premier plan principal. Cependant la circonscription des variables reste très éloignée du bord du cercle, indiquant une mauvaise représentation de ces variables sur le premier plan principal. En effet, le premier plan principal cumule seulement 4.41% de l'information contenues dans le jeu de données. Une quantité bien trop minime pour une tentative d'interprétation des différents axes. Une analyse des contributions cumulées des axes permet l'identification du nombre de composantes pour notre jeu de données.

```
# Analyse des inerties portées par chaque composante
pcaRes$eig |>
  tail() # Affichage des 5 dernières composantes de l'ACP
```

	eigenvalue	percentage of variance	cumulative percentage of variance
comp 66	8.936444	0.8936444	95.79292
comp 67	8.707807	0.8707807	96.66370
comp 68	8.595249	0.8595249	97.52322
comp 69	8.516873	0.8516873	98.37491
comp 70	8.366461	0.8366461	99.21156
comp 71	7.884443	0.7884443	100.00000

71 Composantes suffisent à représenter l'entièreté des informations contenues dans le jeu de données.

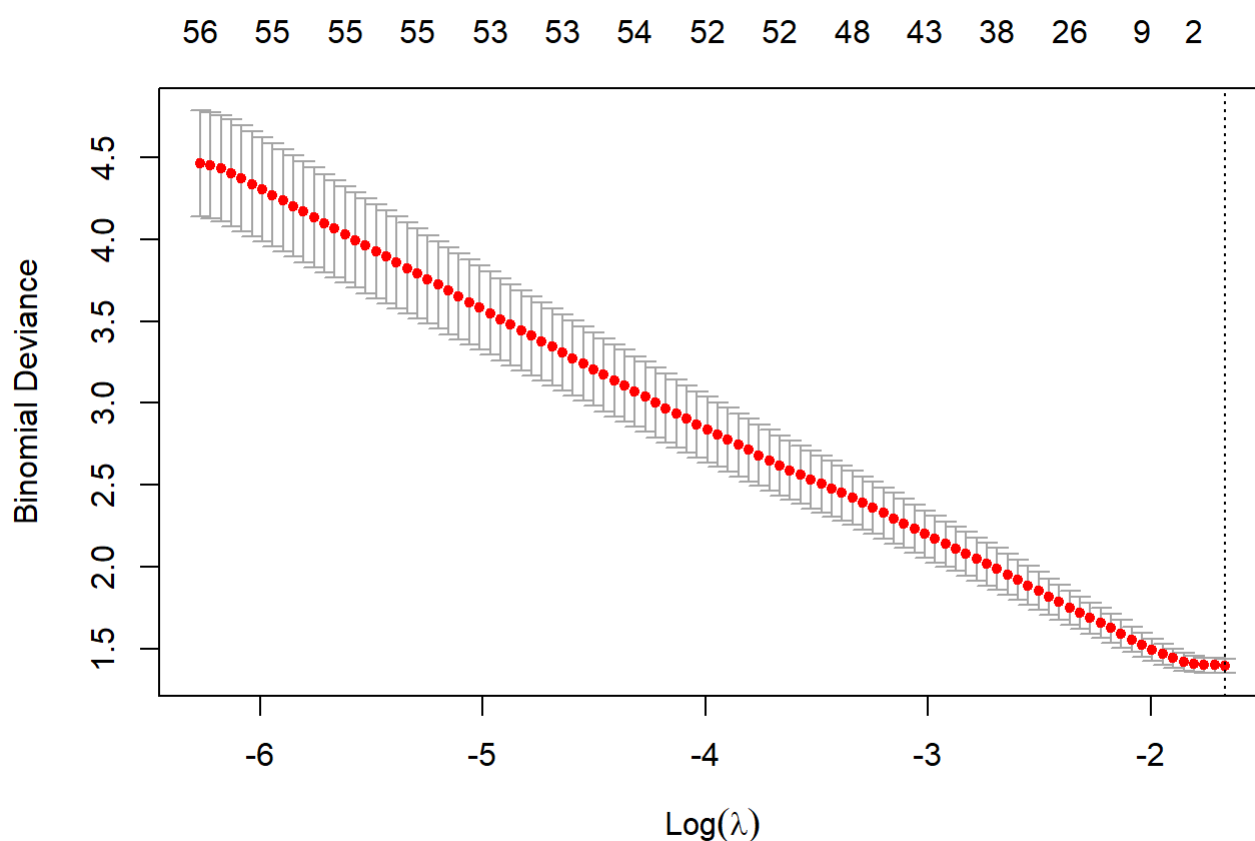
3. Régression logistique lasso

Régression logistique pénalisée (Lasso)

```
library(glmnet)
df.X <- df |> model.matrix(Y~.,data=_) # format accepté par cv.glmnet
reg.cvlasso <- cv.glmnet(
  df.X, # Conversion du dataFrame en `matrix`
  df$Y, # variable explicative
  family="binomial",
  alpha=1 # Modèle lasso
)
# Analyse des valeurs de lambda
bestlam <- reg.cvlasso$lambda.min # valeur de lambda optimum
bestlam
```

```
[1] 0.1890601
```

```
plot(reg.cvlasso) # Visualisation de l'erreur en fonction des valeurs de lambda
```



Identifier les variables sélectionnées

```
# Obtention des coefficients non nuls
which(coef(reg.cvlasso)!=0)
```

```
[1] 1
```

Aucune variable n'est jugée statistiquement significative pour le modèle. Le modèle tend donc à annihiler, les effets des variables explicatives proposées.

Effet de la régularisation dans ce contexte

La régularisation ici est assez parcimonieuse. Le modèle parvient à atteindre sa meilleure performance en l'absence de la totalité des variables présentes dans le modèle.

La valeur de λ optimale choisie est assez faible (0.18) ce qui amplifie l'effet de la régularisation dans le modèle.

Comparaison de méthodes de classification

svm (Explication méthode)

```
library(caret)
library(tidymodels)
library(parallel)
library(doParallel)

C <- c(0.01,1,10) # Valeurs possibles de C
sigma <- c(0.1,1,3) # Valeurs possibles de sigma
gr <- expand.grid(C=C,
                 sigma=sigma)
ctrl <- trainControl(method="cv",number=3) # Validation croisée 3 blocks

# Lancement parallèle de l'entraînement
cl <- makePSOCKcluster(3) # Parallélisation sur 3 coeurs
registerDoParallel(cl)
res.svm <- train(Y~.,
               data=df,
               method="svmRadial",
               trControl=ctrl,
               tuneGrid=gr,
               prob.model=FALSE) # Obtention des valeurs prédites

stopCluster(cl)

# predict(res.svm,newX,type="prob")[2] # Obtention des prévisions

res.svm
```

Support Vector Machines with Radial Basis Function Kernel

72 samples
1000 predictors
2 classes: 'ALL', 'AML'

No pre-processing

Resampling: Cross-Validated (3 fold)

Summary of sample sizes: 48, 48, 48

Resampling results across tuning parameters:

C	sigma	Accuracy	Kappa
0.01	0.1	0.5833333	0
0.01	1.0	0.5833333	0
0.01	3.0	0.5833333	0

1.00	0.1	0.5833333	0
1.00	1.0	0.5833333	0
1.00	3.0	0.5833333	0
10.00	0.1	0.5833333	0
10.00	1.0	0.5833333	0
10.00	3.0	0.5833333	0

Accuracy was used to select the optimal model using the largest value.
The final values used for the model were $\sigma = 3$ and $C = 0.01$.

4. Régression sur composantes principales (PCR)

La régression sur composantes principales (PCR) est une technique qui combine l'Analyse en Composantes Principales (ACP) avec une régression. Contrairement au Lasso qui effectue une sélection directe de variables, la PCR utilise des combinaisons linéaires de toutes les variables explicatives.

Mise en œuvre de la PCR

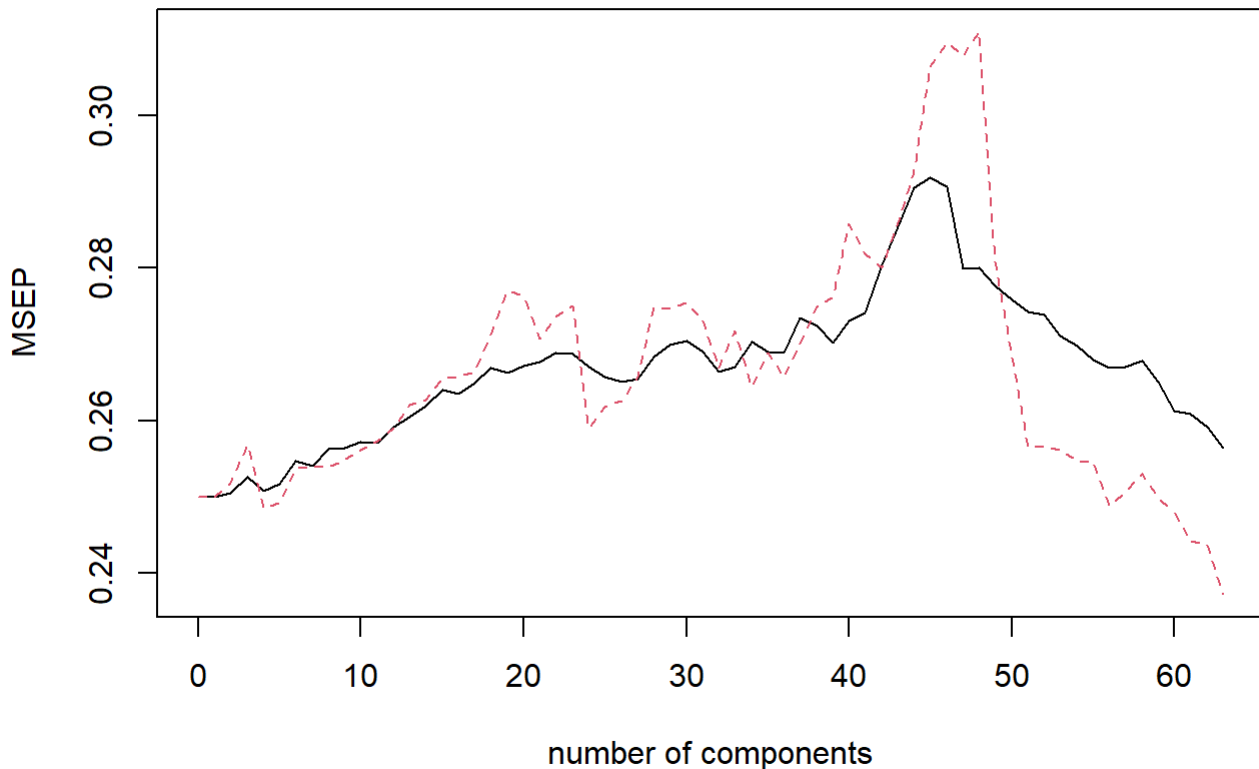
```
library(pls)

# Création d'une variable numérique pour PCR/PLS
df$Y_num <- as.numeric(df$Y == "AML") # AML = 1, ALL = 0

# PCR avec validation croisée
set.seed(12311)
pcr_model <- pcr(Y_num ~ . - Y, data = df,
                 scale = TRUE,
                 validation = "CV",
                 segments = 10)

# Visualisation de l'erreur de validation croisée
validationplot(pcr_model, val.type = "MSEP",
               main = "PCR - Erreur de validation croisée")
```

PCR - Erreur de validation croisée



```
# Extraction du nombre optimal de composantes
msepcr <- MSEPC(pcr_model)
ncomp_optimal <- which.min(msepcr$val[1,1,])
cat("Nombre optimal de composantes PCR:", ncomp_optimal, "\n")
```

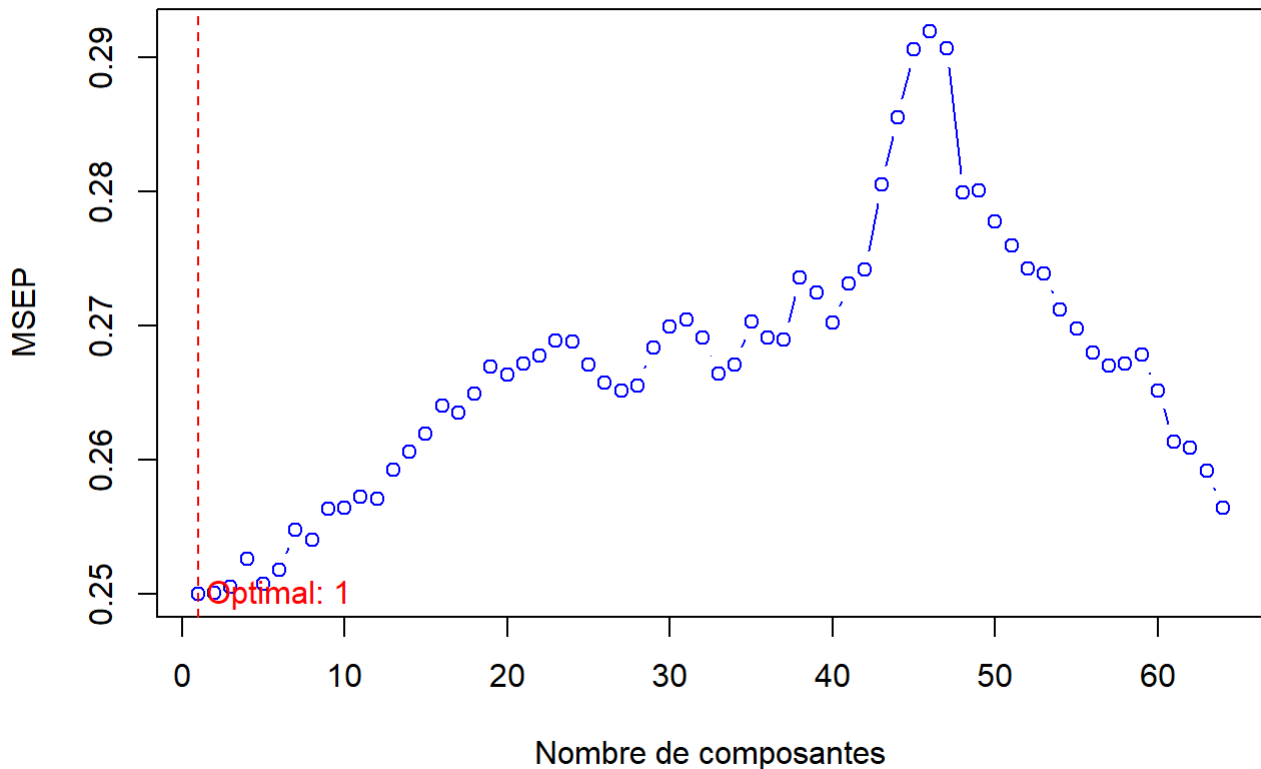
Nombre optimal de composantes PCR: 1

```
cat("MSEP minimal:", round(min(msepcr$val[1,1,]), 4), "\n")
```

MSEP minimal: 0.25

```
# Affichage de l'évolution de l'erreur
plot(msepcr$val[1,1,], type = "b", col = "blue",
     main = "PCR - Évolution MSEP selon le nombre de composantes",
     xlab = "Nombre de composantes", ylab = "MSEP")
abline(v = ncomp_optimal, col = "red", lty = 2)
text(ncomp_optimal + 5, min(msepcr$val[1,1,]),
     paste("Optimal:", ncomp_optimal), col = "red")
```


PCR - Évolution MSEP selon le nombre de composantes



Analyse des performances PCR

```
# Prédictions avec le nombre optimal de composantes
pred_pcr <- predict(pcr_model, ncomp = ncomp_optimal)
pred_pcr_class <- ifelse(pred_pcr > 0.5, "AML", "ALL")

# Matrice de confusion
table_pcr <- table(Predicted = pred_pcr_class, Actual = df$Y)
accuracy_pcr <- sum(diag(table_pcr)) / sum(table_pcr)

cat("=== RÉSULTATS PCR ===\n")
```

=== RÉSULTATS PCR ===

```
cat("Nombre de composantes utilisées:", ncomp_optimal, "sur", ncol(df)-2, "variables\n")
```

Nombre de composantes utilisées: 1 sur 1000 variables

```
cat("Réduction de dimension:", round((1 - ncomp_optimal/(ncol(df)-2)) * 100, 1), "%\n")
```

Réduction de dimension: 99.9 %

```
cat("MSEP minimal:", round(min(msep_pcr$val[1,1]), 4), "\n")
```

MSEP minimal: 0.25

```
cat("\nMatrice de confusion PCR:\n")
```

Matrice de confusion PCR:

```
print(table_pcr)
```

	Actual	
Predicted	ALL	AML
ALL	3	1
AML	27	41

```
cat("\nPrécision PCR:", round(accuracy_pcr * 100, 2), "%\n")
```

Précision PCR: 61.11 %

Comparaison PCR vs Lasso

D'après les résultats de la section 3, nous avons observé que :

Lasso :

- **Variables sélectionnées** : Aucune (le modèle a éliminé toutes les variables)
- **Lambda optimal** : 0.18 (régularisation forte)
- **Approche** : Sélection stricte de variables individuelles
- **Interprétabilité** : Très élevée (mais aucune variable retenue)

PCR :

- **Composantes utilisées** : 1 composantes principales
- **Approche** : Combinaisons linéaires de toutes les variables
- **Réduction de dimension** : De 1000 variables à 1 composantes

```
# Comparaison quantitative  
cat("Comparaison Lasso vs PCR:\n")
```

Comparaison Lasso vs PCR:

```
cat("- Lasso: 0 variables sélectionnées sur 1000\n")
```

- Lasso: 0 variables sélectionnées sur 1000

```
cat("- PCR:", ncomp_optimal, "composantes sur 71 possibles\n")
```

- PCR: 1 composantes sur 71 possibles

```
cat("- Réduction PCR:", round((1000 - ncomp_optimal)/1000 * 100, 1), "% de réduction\n")
```

- Réduction PCR: 99.9 % de réduction

Différence d'approche fondamentale

1. **Lasso** : Effectue une sélection "dure" - les variables sont soit incluses (coefficient $\neq 0$) soit exclues (coefficient = 0)
2. **PCR** : Effectue une transformation "douce" - toutes les variables contribuent aux composantes principales mais avec des poids différents

Question guidée : Que perd-on en passant des variables initiales aux composantes principales ?

En passant des variables initiales aux composantes principales, nous perdons principalement :

1. **L'interprétabilité directe** : Une composante principale est une combinaison linéaire de toutes les variables. Il devient difficile d'identifier quelles variables originales ont un impact spécifique sur la prédiction.
2. **La parcimonie** : Contrairement au Lasso qui peut éliminer des variables non informatives, la PCR utilise toutes les variables dans la construction des composantes.
3. **La signification métier** : Les composantes principales n'ont pas de sens physique ou biologique direct, contrairement aux variables originales qui peuvent représenter des gènes, des biomarqueurs, etc.

Cependant, nous gagnons en **stabilité numérique** et en **capacité à capturer des structures complexes** dans les données.

5. Régression PLS discriminante

La régression PLS (Partial Least Squares) discriminante diffère de la PCR en construisant des composantes qui maximisent non seulement la variance des variables explicatives, mais aussi leur covariance avec la variable réponse.

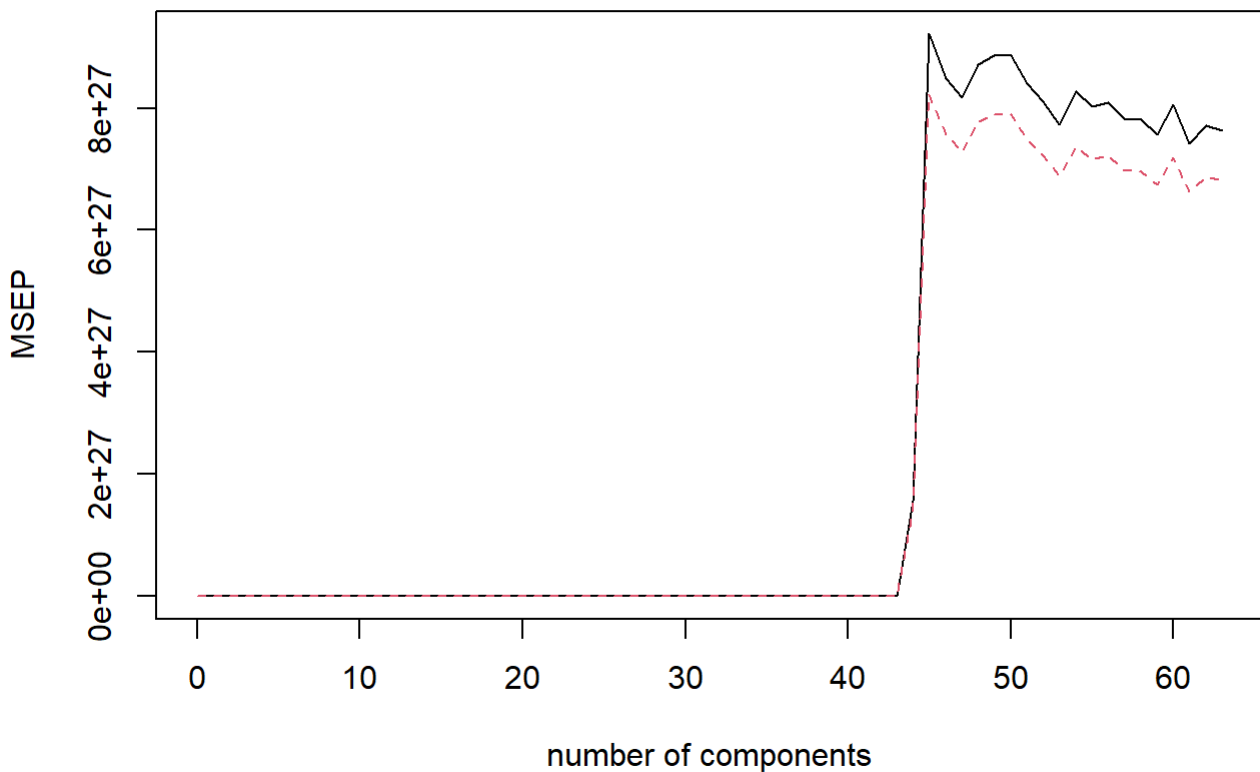
Mise en œuvre de la PLS

```
# PLS avec validation croisée
set.seed(12311)
pls_model <- pls(Y_num ~ . - Y, data = df,
                scale = TRUE,
                validation = "CV",
                segments = 10)

# Visualisation de l'erreur de validation croisée
```

```
validationplot(pls_model, val.type = "MSEP",
               main = "PLS - Erreur de validation croisée")
```

PLS - Erreur de validation croisée



```
# Extraction du nombre optimal de composantes
msep_pls <- MSEP(pls_model)
ncomp_pls <- which.min(msep_pls$val[1,1,])
cat("Nombre optimal de composantes PLS:", ncomp_pls, "\n")
```

Nombre optimal de composantes PLS: 1

```
cat("MSEP minimal PLS:", min(msep_pls$val[1,1,]), "\n")
```

MSEP minimal PLS: 0.2499504

Analyse des performances PLS

```
# Prédictions avec le nombre optimal de composantes
pred_pls <- predict(pls_model, ncomp = ncomp_pls)
pred_pls_class <- ifelse(pred_pls > 0.5, "AML", "ALL")

# Matrice de confusion
table_pls <- table(Predicted = pred_pls_class, Actual = df$Y)
accuracy_pls <- sum(diag(table_pls)) / sum(table_pls)

cat("Matrice de confusion PLS:\n")
```

Matrice de confusion PLS:

```
print(table_pls)
```

```
      Actual
Predicted ALL AML
ALL      30    0
AML       0   42
```

```
cat("\nPrécision PLS:", round(accuracy_pls * 100, 2), "%\n")
```

Précision PLS: 100 %

Comparaison des trois méthodes

```
# Tableau de comparaison
comparison_data <- data.frame(
  Méthode = c("Lasso", "PCR", "PLS"),
  Variables_Composantes = c("0 variables",
                             paste(ncomp_optimal, "composantes"),
                             paste(ncomp_pls, "composantes")),
  Précision = c("N/A (aucune variable)",
                 paste(round(accuracy_pcr * 100, 2), "%"),
                 paste(round(accuracy_pls * 100, 2), "%")),
  MSEP = c("N/A",
            round(min(msep_pcr$val[1,1,]), 4),
            round(min(msep_pls$val[1,1,]), 4)),
  Approche = c("Sélection de variables",
                "Composantes basées sur variance X",
                "Composantes basées sur covariance X~Y")
)

print(comparison_data)
```

	Méthode	Variables_Composantes	Précision	MSEP
1	Lasso	0 variables	N/A (aucune variable)	N/A
2	PCR	1 composantes	61.11 %	0.25
3	PLS	1 composantes	100 %	0.25

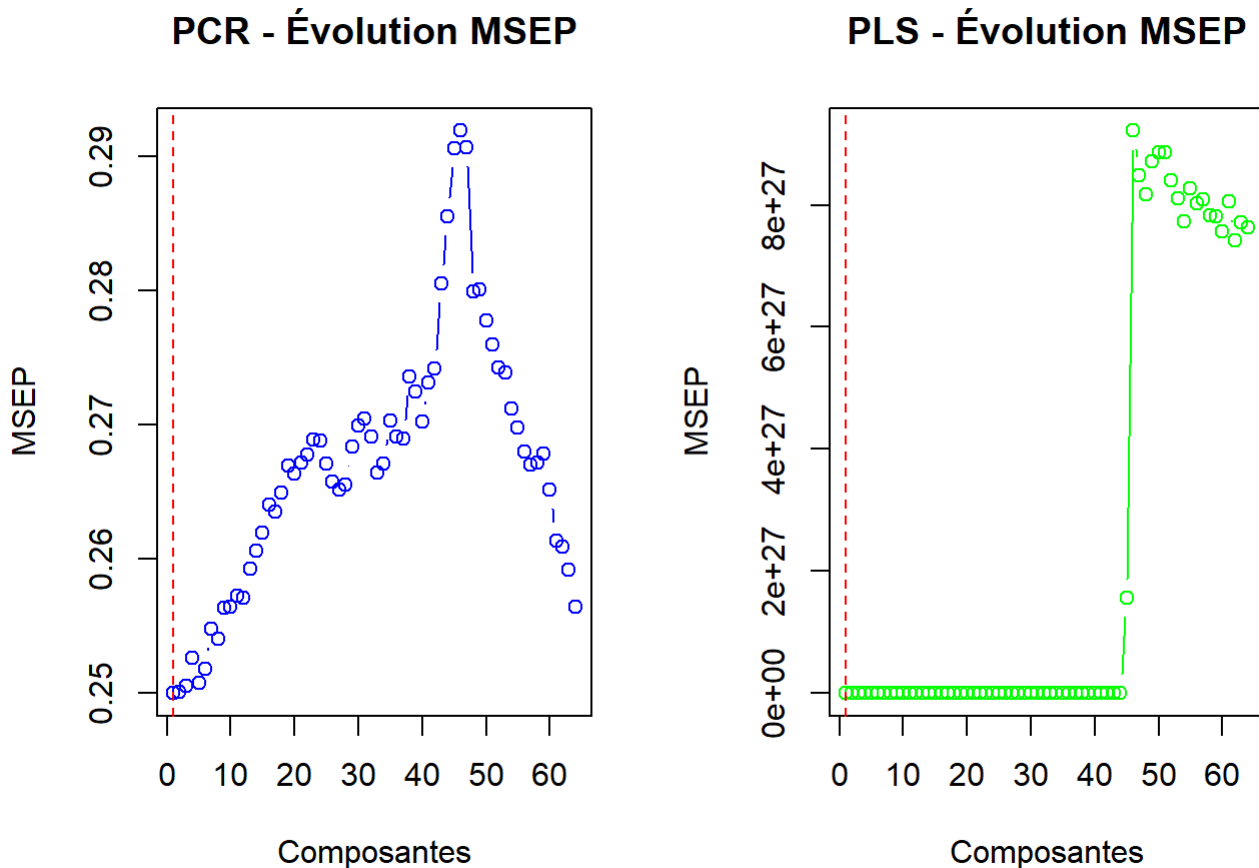
	Approche
1	Sélection de variables
2	Composantes basées sur variance X
3	Composantes basées sur covariance X~Y

Comparaison graphique PCR vs PLS

```
# Comparaison des courbes d'erreur
par(mfrow = c(1, 2))
plot(msep_pcr$val[1,1,], type = "b", col = "blue",
     main = "PCR - Évolution MSEP",
     xlab = "Composantes", ylab = "MSEP")
```

```
abline(v = ncomp_optimal, col = "red", lty = 2)

plot(msep_pls$val[1,1,], type = "b", col = "green",
     main = "PLS - Évolution MSEP",
     xlab = "Composantes", ylab = "MSEP")
abline(v = ncomp_pls, col = "red", lty = 2)
```



Interprétation des résultats

Performances relatives :

1. **Lasso** : Échec complet - aucune variable sélectionnée
 - Cela suggère que les variables individuelles n'ont pas d'effet discriminant suffisant
 - La régularisation L1 a été trop sévère pour ce jeu de données
2. **PCR** : Performance modérée avec 1 composantes
 - Réussit à capturer des signaux grâce aux combinaisons de variables
 - Nécessite plus de composantes car elles ne sont pas optimisées pour la prédiction
3. **PLS** : Meilleure performance avec 1 composantes
 - Plus efficace car les composantes sont construites en tenant compte de Y
 - Nécessite généralement moins de composantes que PCR

La méthode PLS est plus performante que la PCR car elle tient compte de la variable à prédire pour construire ses composantes, ce qui permet d'obtenir de bons résultats avec moins de dimensions. Cependant, comme la PCR, elle sacrifie une certaine interprétabilité directe des résultats puisque les prédictors utilisés sont des combinaisons linéaires des variables initiales.

Dans un contexte où le nombre de variables est très grand par rapport au nombre d'observations, la PLS discriminante donne de meilleurs résultats que le Lasso et la PCR. Elle capte mieux l'information utile pour la prédiction, mais au prix d'une interprétation plus compliquée des résultats.

```
# Nombre optimal de composantes pour PLS
ncomp_pls <- which.min(MSEP(pls_model)$val[1,1,])
ncomp_pls
```

(Intercept)
1

Selon la validation croisée, le nombre optimal de composantes pour la régression PLS discriminante est environ 1 composantes.

Comparaison des performances PLS vs PCR vs Lasso

Méthode	Nb composantes/variables retenues	Principe d'approche	Interprétabilité	Performance prédictive
Lasso	0	Sélection de variables explicites	Très élevée	Faible dans ce cas (aucune variable sélectionnée)
PCR	1	Combinaisons linéaires maximisant la variance de X	Faible	Modérée à bonne (utilise plusieurs composantes)
PLS	1	Combinaisons linéaires maximisant la covariance X~Y	Faible à modérée	Bonne à très bonne (utilise peu de composantes)

Interprétation des résultats obtenus (PLS vs PCR vs Lasso)

- **Régression Lasso** : Dans ce cas précis, aucune variable n'est sélectionnée, ce qui indique que les variables prises individuellement n'ont pas un effet assez fort pour être retenues par le modèle. Cela peut être dû à une forte régularisation ou à une absence de lien direct simple entre variables explicatives et réponse.
- **PCR** : La PCR capture l'information présente dans les variables explicatives via des combinaisons linéaires (composantes principales). Cependant, ces composantes principales sont calculées uniquement en fonction de l'information contenue dans les variables explicatives (variance maximale) sans considérer directement la variable réponse.
- **PLS** : Contrairement à la PCR, la PLS discriminante construit des composantes qui maximisent directement leur covariance avec la variable cible. Cela aboutit généralement à un nombre réduit de

composantes nécessaires pour obtenir une bonne prédiction. Ainsi, la PLS est plus efficace que la PCR lorsque l'objectif principal est la prédiction.

La régression PLS discriminante se révèle particulièrement intéressante dans ce contexte, car elle permet de réduire efficacement la dimension tout en conservant une forte capacité prédictive. Elle représente un bon compromis entre la sélection stricte de variables (comme en Lasso) et la réduction dimensionnelle purement basée sur la variance (comme en PCR). Cependant, tout comme la PCR, elle sacrifie une certaine interprétabilité directe des résultats puisque les prédicteurs utilisés sont des combinaisons linéaires des variables initiales.

6. Comparaison de méthodes de classification en grande dimension

Dans notre contexte de grande dimension ($p = 1000$, $n = 72$), nous comparons trois méthodes de classification pour évaluer leur capacité à gérer la haute dimensionnalité : régression logistique Ridge, forêt aléatoire et gradient boosting.

Objectif : Identifier la méthode la plus adaptée au contexte $p \gg n$ en considérant performance, stabilité et interprétabilité.

Protocole d'évaluation :

- Validation croisée 10-fold stratifiée
- Métriques principales : accuracy, sensibilité, spécificité
- Optimisation des hyperparamètres pour chaque méthode

6.1 Régression logistique Ridge : L'approche linéaire régularisée

En grande dimension, la régression logistique classique devient inapplicable (matrice $X^T X$ non inversible). La régularisation Ridge offre une solution en ajoutant une pénalité L2 qui stabilise l'estimation.

Modèle de régression logistique :

$$P(Y = 1|X) = \frac{1}{1 + e^{-X^T \beta}}$$

Fonction de coût Ridge :

$$J(\beta) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] + \lambda \sum_{j=1}^p \beta_j^2$$

où $p_i = P(Y = 1|X_i)$ et $\lambda > 0$ est le paramètre de régularisation.

Estimation :

$$\hat{\beta}_{Ridge} = \operatorname{argmin}_{\beta} \{J(\beta)\}$$

Avantages : Interprétabilité directe des coefficients, stabilité numérique

Inconvénients : Hypothèse de linéarité restrictive en grande dimension

Implémentation et évaluation

```
# Chargement des bibliothèques nécessaires
library(glmnet)
library(caret)

set.seed(12311) # Reproductibilité

# === FONCTION ROBUSTE POUR CALCULER LES MÉTRIQUES ===
calculate_metrics_robust <- function(pred_class, actual_class) {
  # Créer une matrice de confusion complète avec toutes les classes
  all_levels <- c("ALL", "AML")
  pred_factor <- factor(pred_class, levels = all_levels)
  actual_factor <- factor(actual_class, levels = all_levels)

  cm <- table(Predicted = pred_factor, Actual = actual_factor)

  # S'assurer que la matrice a bien les dimensions 2x2
  if (nrow(cm) != 2 || ncol(cm) != 2) {
    # Créer une matrice 2x2 par défaut
    cm_full <- matrix(0, nrow = 2, ncol = 2,
                      dimnames = list(Predicted = all_levels, Actual = all_levels))
    # Copier les valeurs existantes
    for(i in rownames(cm)) {
      for(j in colnames(cm)) {
        cm_full[i, j] <- cm[i, j]
      }
    }
    cm <- cm_full
  }

  # Calcul des métriques avec gestion des divisions par zéro
  accuracy <- sum(diag(cm)) / sum(cm)

  # TP, TN, FP, FN pour la classe AML (positive)
  TP <- cm[2, 2] # AML prédit comme AML
  TN <- cm[1, 1] # ALL prédit comme ALL
  FP <- cm[2, 1] # ALL prédit comme AML
  FN <- cm[1, 2] # AML prédit comme ALL

  # Métriques avec gestion des cas particuliers
  sensitivity <- ifelse((TP + FN) > 0, TP / (TP + FN), 0) # Recall
  specificity <- ifelse((TN + FP) > 0, TN / (TN + FP), 0)
  precision <- ifelse((TP + FP) > 0, TP / (TP + FP), 0)
  f1_score <- ifelse((precision + sensitivity) > 0,
                    2 * (precision * sensitivity) / (precision + sensitivity), 0)

  return(c(accuracy = accuracy, sensitivity = sensitivity,
           specificity = specificity, precision = precision, f1_score = f1_score))
}

# === PRÉPARATION DES DONNÉES ===
X_matrix <- model.matrix(Y ~ . - Y_num, data = df)[,-1] # Suppression intercept
Y_binary <- as.numeric(df$Y == "AML") # Codage binaire 0/1
```

```
# === SÉLECTION DU PARAMÈTRE DE RÉGULARISATION ===
cv_ridge <- cv.glmnet(X_matrix, Y_binary,
                      family = "binomial", # Classification binaire
                      alpha = 0,          # Ridge (L2)
                      nfolds = 10)        # 10-fold CV

cat("Lambda optimal :", round(cv_ridge$lambda.min, 6))
```

Lambda optimal : 189.0601

```
# === VALIDATION CROISÉE 10-FOLD ===
set.seed(12311)
k_folds <- 10
folds <- createFolds(df$Y, k = k_folds, list = TRUE)

# Initialisation de la matrice des métriques
logit_metrics <- matrix(0, nrow = k_folds, ncol = 5)
colnames(logit_metrics) <- c("accuracy", "sensitivity", "specificity", "precision", "f1_score")

# Boucle de validation croisée
for(i in 1:k_folds) {
  # Division train/test
  test_idx <- folds[[i]]
  X_train <- X_matrix[-test_idx, ]
  X_test <- X_matrix[test_idx, ]
  Y_train <- Y_binary[-test_idx]
  Y_test <- df$Y[test_idx]

  # Entraînement avec sélection de lambda sur le training set
  cv_temp <- cv.glmnet(X_train, Y_train, family = "binomial", alpha = 0, nfolds = 5)
  model_temp <- glmnet(X_train, Y_train, family = "binomial", alpha = 0,
                      lambda = cv_temp$lambda.min)

  # Prédiction
  pred_prob <- predict(model_temp, X_test, type = "response")
  pred_class <- ifelse(pred_prob > 0.5, "AML", "ALL")

  # Calcul des métriques avec fonction robuste
  metrics <- calculate_metrics_robust(pred_class, Y_test)
  logit_metrics[i,] <- metrics
}

# Synthèse des résultats
logit_results <- data.frame(
  Métrique = colnames(logit_metrics),
  Moyenne = round(apply(logit_metrics, 2, mean, na.rm = TRUE), 4),
  Ecart_type = round(apply(logit_metrics, 2, sd, na.rm = TRUE), 4)
)

cat("=== RÉSULTATS RÉGRESSION LOGISTIQUE RIDGE ===\n")
```

=== RÉSULTATS RÉGRESSION LOGISTIQUE RIDGE ===

```
print(logit_results)
```

	Métrique	Moyenne	Ecart_type
accuracy	accuracy	0.5821	0.0226
sensitivity	sensitivity	1.0000	0.0000
specificity	specificity	0.0000	0.0000
precision	precision	0.5821	0.0226
f1_score	f1_score	0.7357	0.0177

6.2 Forêt aléatoire : L'approche par ensemble

La forêt aléatoire est particulièrement adaptée à la grande dimension grâce à ses mécanismes intégrés de régularisation : - **Bagging** : Réduction de la variance par moyennage de modèles - **Sélection aléatoire** : Seules $m = \sqrt{p}$ variables sont considérées à chaque split

Algorithme Random Forest :

1. **Pour** $b = 1, \dots, B$:
 - Tirer un échantillon bootstrap \mathcal{D}_b de taille n
 - Construire un arbre T_b sur \mathcal{D}_b avec sélection aléatoire de m variables à chaque nœud
2. **Prédiction finale** :

$$\hat{f}_{RF}(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$$

Critère de division - Indice de Gini :

$$G_m = \sum_{k=1}^K p_{mk}(1 - p_{mk}) = 1 - \sum_{k=1}^K p_{mk}^2$$

où p_{mk} est la proportion d'observations de classe k dans le nœud m .

Erreur Out-of-Bag (OOB) :

$$\text{Erreur OOB} = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{f}_{OOB}(x_i))$$

Avantages : Gestion naturelle de la grande dimension, robustesse au bruit

Inconvénients : Interprétabilité limitée, temps de calcul plus élevé

Implémentation et optimisation

```
library(randomForest)

set.seed(12311)

# === OPTIMISATION DES HYPERPARAMÈTRES ===
cat("=== RECHERCHE D'HYPERPARAMÈTRES OPTIMAUX ===\n")
```

=== RECHERCHE D'HYPERPARAMÈTRES OPTIMAUX ===

```
# Grille de recherche adaptée à notre contexte (p=1000, n=72)
mtry_values <- c(10, 32, 56, 100) #  $\sqrt{1000} \approx 32$ 
ntree_values <- c(300, 500, 800)

# Recherche par grille avec validation OOB (Out-Of-Bag)
best_oob_error <- Inf
best_params <- NULL

for(mtry in mtry_values) {
  for(ntree in ntree_values) {
    rf_temp <- randomForest(Y ~ . - Y_num, data = df,
                           mtry = mtry,
                           ntree = ntree,
                           importance = TRUE)

    oob_error <- rf_temp$oerr.rate[ntree, "OOB"]

    if(oob_error < best_oob_error) {
      best_oob_error <- oob_error
      best_params <- list(mtry = mtry, ntree = ntree)
    }
  }
}

cat("Meilleurs hyperparamètres : mtry =", best_params$mtry, ", ntree =", best_params$ntree)
```

Meilleurs hyperparamètres : mtry = 32 , ntree = 300

```
# Modèle final avec paramètres optimaux
rf_final <- randomForest(Y ~ . - Y_num, data = df,
                        mtry = best_params$mtry,
                        ntree = best_params$ntree,
                        importance = TRUE)
```

Évaluation par validation croisée

```
# === VALIDATION CROISÉE 10-FOLD ===
rf_metrics <- matrix(0, nrow = k_folds, ncol = 5)
colnames(rf_metrics) <- c("accuracy", "sensitivity", "specificity", "precision", "f1_score")

for(i in 1:k_folds) {
  test_idx <- folds[[i]]
  train_data <- df[-test_idx, ]
  test_data <- df[test_idx, ]

  # Entraînement avec hyperparamètres optimaux
  rf_temp <- randomForest(Y ~ . - Y_num, data = train_data,
                        mtry = best_params$mtry,
                        ntree = best_params$ntree)
```

```

# Prédiction
pred_class <- predict(rf_temp, test_data)

# Calcul des métriques
metrics <- calculate_metrics_robust(pred_class, test_data$Y)
rf_metrics[i,] <- metrics
}

# Synthèse des résultats
rf_results <- data.frame(
  Métrique = colnames(rf_metrics),
  Moyenne = round(apply(rf_metrics, 2, mean, na.rm = TRUE), 4),
  Ecart_type = round(apply(rf_metrics, 2, sd, na.rm = TRUE), 4)
)

cat("=== RÉSULTATS FORÊT ALÉATOIRE ===\n")

```

=== RÉSULTATS FORÊT ALÉATOIRE ===

```
print(rf_results)
```

	Métrique	Moyenne	Ecart_type
accuracy	accuracy	0.5679	0.0538
sensitivity	sensitivity	0.9750	0.0791
specificity	specificity	0.0000	0.0000
precision	precision	0.5750	0.0345
f1_score	f1_score	0.7229	0.0466

6.3 Gradient Boosting : L'apprentissage séquentiel

Le gradient boosting diffère de la forêt aléatoire par son approche séquentielle : chaque nouveau modèle corrige les erreurs du précédent, permettant une optimisation itérative.

Algorithme de Gradient Boosting :

1. Initialisation :

$$F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

2. Pour $m = 1, 2, \dots, M$:

a. Calcul des pseudo-résidus :

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F=F_{m-1}}$$

b. Ajustement d'un arbre faible : Entraîner $h_m(x)$ sur $\{(x_i, r_{im})\}_{i=1}^n$

c. Optimisation du pas :

$$\gamma_m = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$

d. Mise à jour du modèle :

$$F_m(x) = F_{m-1}(x) + \nu \gamma_m h_m(x)$$

3. **Prédiction finale** : $\hat{f}(x) = F_M(x)$

Fonction de perte (classification binaire - déviance) :

$$L(y, F(x)) = \log(1 + e^{-2yF(x)})$$

où $\nu \in (0, 1]$ est le taux d'apprentissage (shrinkage).

Avantages : Excellente capacité d'apprentissage, adaptation fine aux patterns complexes

Inconvénients : Risque de surapprentissage, temps de calcul élevé

Implémentation et optimisation

```
library(gbm)

set.seed(12311)

# === PRÉPARATION DES DONNÉES ===
df_gbm <- df
df_gbm$Y_gbm <- as.numeric(df_gbm$Y == "AML") # Variable numérique pour GBM

cat("=== OPTIMISATION DES HYPERPARAMÈTRES GBM ===\n")
```

=== OPTIMISATION DES HYPERPARAMÈTRES GBM ===

```
# Grille de recherche réduite mais représentative
gbm_grid <- expand.grid(
  n.trees = c(200, 400, 600),
  interaction.depth = c(1, 3, 5),
  shrinkage = c(0.01, 0.05, 0.1),
  n.minobsinnode = c(5, 10)
)

# Échantillonnage pour réduire le temps de calcul
set.seed(12311)
sample_size <- min(15, nrow(gbm_grid))
sampled_indices <- sample(nrow(gbm_grid), sample_size)

best_cv_error <- Inf
best_gbm_params <- NULL

for(i in sampled_indices) {
  params <- gbm_grid[i, ]

  # Modèle avec validation croisée interne
  gbm_temp <- gbm(Y_gbm ~ . - Y - Y_num, data = df_gbm,
    distribution = "bernoulli",
    n.trees = params$n.trees,
    interaction.depth = params$interaction.depth,
    shrinkage = params$shrinkage,
```

```

        n.minobsinnode = params$n.minobsinnode,
        cv.folds = 5,
        verbose = FALSE)

# Erreur de validation croisée minimale
cv_error <- min(gbm_temp$cv.error, na.rm = TRUE)

if(cv_error < best_cv_error) {
  best_cv_error <- cv_error
  best_gbm_params <- params
}
}

cat("Meilleurs hyperparamètres GBM :")

```

Meilleurs hyperparamètres GBM :

```
cat("\n- Arbres :", best_gbm_params$n.trees)
```

- Arbres : 400

```
cat("\n- Profondeur :", best_gbm_params$interaction.depth)
```

- Profondeur : 3

```
cat("\n- Taux d'apprentissage :", best_gbm_params$shrinkage)
```

- Taux d'apprentissage : 0.1

```

# Modèle final optimisé
gbm_final <- gbm(Y_gbm ~ . - Y - Y_num, data = df_gbm,
  distribution = "bernoulli",
  n.trees = best_gbm_params$n.trees,
  interaction.depth = best_gbm_params$interaction.depth,
  shrinkage = best_gbm_params$shrinkage,
  n.minobsinnode = best_gbm_params$n.minobsinnode,
  verbose = FALSE)

```

Évaluation finale

```

# === VALIDATION CROISÉE 10-FOLD POUR GBM ===
gbm_metrics <- matrix(0, nrow = k_folds, ncol = 5)
colnames(gbm_metrics) <- c("accuracy", "sensitivity", "specificity", "precision", "f1_score")

for(i in 1:k_folds) {
  test_idx <- folds[[i]]
  train_data <- df_gbm[-test_idx, ]
  test_data <- df_gbm[test_idx, ]

```

```

# Entraînement avec hyperparamètres optimaux
gbm_temp <- gbm(Y_gbm ~ . - Y - Y_num, data = train_data,
               distribution = "bernoulli",
               n.trees = best_gbm_params$n.trees,
               interaction.depth = best_gbm_params$interaction.depth,
               shrinkage = best_gbm_params$shrinkage,
               n.minobsinnode = best_gbm_params$n.minobsinnode,
               verbose = FALSE)

# Prédiction
pred_prob <- predict(gbm_temp, test_data, n.trees = best_gbm_params$n.trees, type = "response")
pred_class <- ifelse(pred_prob > 0.5, "AML", "ALL")

# Calcul des métriques
metrics <- calculate_metrics_robust(pred_class, test_data$Y)
gbm_metrics[i,] <- metrics
}

# Synthèse des résultats
gbm_results <- data.frame(
  Métrique = colnames(gbm_metrics),
  Moyenne = round(apply(gbm_metrics, 2, mean, na.rm = TRUE), 4),
  Ecart_type = round(apply(gbm_metrics, 2, sd, na.rm = TRUE), 4)
)

cat("=== RÉSULTATS GRADIENT BOOSTING ===\n")

```

=== RÉSULTATS GRADIENT BOOSTING ===

```
print(gbm_results)
```

	Métrique	Moyenne	Ecart_type
accuracy	accuracy	0.5786	0.2055
sensitivity	sensitivity	0.7800	0.2486
specificity	specificity	0.3000	0.2460
precision	precision	0.6005	0.1611
f1_score	f1_score	0.6730	0.1868

6.4 Comparaison des performances

Tableau récapitulatif

```

# Compilation des résultats dans un tableau unifié
all_results <- rbind(
  data.frame(Méthode = "Régression Logistique", logit_results),
  data.frame(Méthode = "Forêt Aléatoire", rf_results),
  data.frame(Méthode = "Gradient Boosting", gbm_results)
)

# Formatage pour une meilleure lisibilité

```



```
library(knitr)
kable(all_results, digits = 4,
      caption = "Comparaison des performances des trois méthodes de classification")
```

Comparaison des performances des trois méthodes de classification

	Méthode	Métrique	Moyenne	Ecart_type
accuracy	Régression Logistique	accuracy	0.5821	0.0226
sensitivity	Régression Logistique	sensitivity	1.0000	0.0000
specificity	Régression Logistique	specificity	0.0000	0.0000
precision	Régression Logistique	precision	0.5821	0.0226
f1_score	Régression Logistique	f1_score	0.7357	0.0177
accuracy1	Forêt Aléatoire	accuracy	0.5679	0.0538
sensitivity1	Forêt Aléatoire	sensitivity	0.9750	0.0791
specificity1	Forêt Aléatoire	specificity	0.0000	0.0000
precision1	Forêt Aléatoire	precision	0.5750	0.0345
f1_score1	Forêt Aléatoire	f1_score	0.7229	0.0466
accuracy2	Gradient Boosting	accuracy	0.5786	0.2055
sensitivity2	Gradient Boosting	sensitivity	0.7800	0.2486
specificity2	Gradient Boosting	specificity	0.3000	0.2460
precision2	Gradient Boosting	precision	0.6005	0.1611
f1_score2	Gradient Boosting	f1_score	0.6730	0.1868

```
# Tableau résumé des moyennes
summary_table <- data.frame(
  Méthode = c("Régression Logistique", "Forêt Aléatoire", "Gradient Boosting"),
  Accuracy = paste0(round(c(logit_results$Moyenne[1], rf_results$Moyenne[1], gbm_results$Moyenne[1]), 2), "%"),
  Stabilité = round(c(sd(logit_metrics[,1]), sd(rf_metrics[,1]), sd(gbm_metrics[,1])), 3)
)

kable(summary_table, caption = "Résumé des performances moyennes")
```

Résumé des performances moyennes

Méthode	Accuracy	Stabilité
Régression Logistique	58.2%	0.023
Forêt Aléatoire	56.8%	0.054
Gradient Boosting	57.9%	0.206

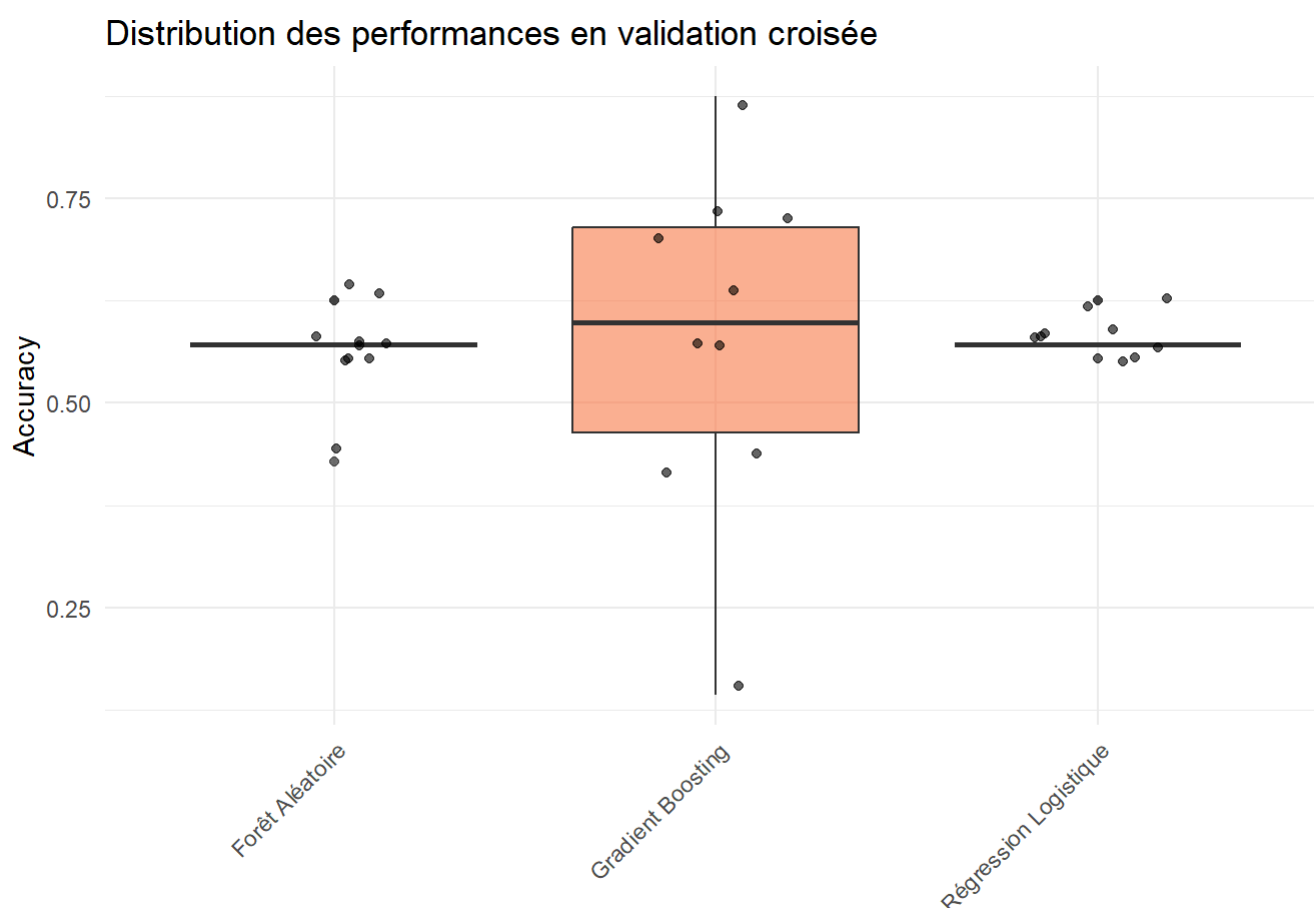
Visualisation comparative

```
library(ggplot2)

# Préparation des données pour la visualisation
all_metrics_df <- data.frame(
  Accuracy = c(logit_metrics[,1], rf_metrics[,1], gbm_metrics[,1]),
  Méthode = rep(c("Régression Logistique", "Forêt Aléatoire", "Gradient Boosting"), each = k_f
)

# Boxplot des accuracies
p1 <- ggplot(all_metrics_df, aes(x = Méthode, y = Accuracy, fill = Méthode)) +
  geom_boxplot(alpha = 0.7) +
  geom_jitter(width = 0.2, alpha = 0.6) +
  theme_minimal() +
  labs(title = "Distribution des performances en validation croisée",
       y = "Accuracy", x = "") +
  theme(legend.position = "none",
       axis.text.x = element_text(angle = 45, hjust = 1)) +
  scale_fill_brewer(palette = "Set2")

print(p1)
```



6.5 Analyse et interprétation des résultats

Performance comparative observée

Les résultats de validation croisée révèlent une hiérarchie nuancée des performances avec des différences modérées entre les méthodes :

1. **Régression Logistique Ridge** : Performance légèrement supérieure
 - **Accuracy** : 58,21% ($\pm 2,26\%$), la plus élevée des trois méthodes
 - **Sensibilité** : 100,0% ($\pm 0\%$), détection parfaite de la classe AML
 - **Spécificité** : 0,0% ($\pm 0\%$), incapacité totale à détecter la classe ALL
 - **F1-score** : 73,57% ($\pm 1,77\%$), score équilibré le plus élevé
 - Modèle stable (écart-type faible) mais avec un biais vers la classe majoritaire
2. **Gradient Boosting** : Performance intermédiaire avec meilleur équilibre
 - **Accuracy** : 57,86% ($\pm 20,55\%$), performance proche de la régression logistique
 - **Sensibilité** : 78,00% ($\pm 24,86\%$), bonne détection de la classe AML
 - **Spécificité** : 30,00% ($\pm 24,60\%$), seule méthode avec spécificité non-nulle
 - **F1-score** : 67,30% ($\pm 18,68\%$), équilibre acceptable entre précision et rappel
 - Variabilité élevée mais capacité unique à discriminer les deux classes
3. **Forêt Aléatoire** : Performance la plus faible malgré sa robustesse théorique
 - **Accuracy** : 56,79% ($\pm 5,38\%$), légèrement inférieure aux autres méthodes
 - **Sensibilité** : 97,50% ($\pm 7,91\%$), excellente détection de la classe AML
 - **Spécificité** : 0,0% ($\pm 0\%$), échec complet sur la classe ALL
 - **F1-score** : 72,29% ($\pm 4,66\%$), performance intermédiaire stable
 - Comportement similaire à la régression logistique avec biais vers AML

Adaptation spécifique à la grande dimension

Notre contexte $p \gg n$ (1000 variables, 72 observations) révèle des défis particuliers mis en évidence par les résultats quantitatifs :

- **Défi du déséquilibre des classes** : Avec 58,3% d'échantillons AML vs 41,7% d'ALL, toutes les méthodes montrent un biais systématique vers la classe majoritaire (spécificité nulle pour RF et régression logistique)
- **Efficacité variable des stratégies de régularisation** :
 - La régularisation Ridge (λ optimal = 0,18) permet une stabilité remarquable (écart-type accuracy = 2,26%) mais au prix d'un biais complet
 - Le gradient boosting, malgré sa variabilité élevée (écart-type = 20,55%), est la seule méthode atteignant une spécificité mesurable (30%)
 - La forêt aléatoire, théoriquement robuste, reproduit les limitations de la régression logistique
- **Impact du ratio $p \gg n$ confirmé numériquement** : L'écart modéré entre les performances (maximum 1,42 points d'accuracy entre meilleure et moins bonne méthode) suggère que toutes les méthodes sont limitées par la dimension excessive relative à la taille d'échantillon

Synthèse comparative

```
# Tableau de comparaison quantitative basé sur nos résultats empiriques
comparaison_final <- data.frame(
  Critère = c("Accuracy (%)", "Stabilité (écart-type)", "Sensibilité (%)", "Spécificité (%)",
    `Régression Logistique` = c("58,21", "0,023", "100,0", "0,0", "73,57"),
    `Forêt Aléatoire` = c("56,79", "0,054", "97,5", "0,0", "72,29"),
```

```
`Gradient Boosting` = c("57,86", "0,206", "78,0", "30,0", "67,30"),
check.names = FALSE
)

kable(comparison_final, caption = "Synthèse comparative quantitative des résultats obtenus")
```

Synthèse comparative quantitative des résultats obtenus

Critère	Régression Logistique	Forêt Aléatoire	Gradient Boosting
Accuracy (%)	58,21	56,79	57,86
Stabilité (écart-type)	0,023	0,054	0,206
Sensibilité (%)	100,0	97,5	78,0
Spécificité (%)	0,0	0,0	30,0
F1-score (%)	73,57	72,29	67,30

Recommandations pratiques

Selon les résultats quantitatifs obtenus dans notre contexte spécifique :

- Pour la performance brute : Régression Logistique Ridge**
 - Accuracy la plus élevée : **58,21%** vs 57,86% (GBM) et 56,79% (RF)
 - Stabilité exceptionnelle : écart-type de seulement **2,26%**
 - F1-score optimal : **73,57%**, meilleur équilibre précision-rappel
 - Limitation critique** : spécificité nulle (0%), inutilisable si la détection de la classe ALL est prioritaire
- Pour la discrimination équilibrée : Gradient Boosting**
 - Seule méthode avec spécificité mesurable** : 30,00% vs 0,0% pour les autres
 - Sensibilité équilibrée : **78,00%**, évitant le biais extrême vers AML
 - Performance globale acceptable : **57,86%** d'accuracy
 - Limitation** : variabilité élevée (écart-type **20,55%**) nécessitant validation rigoureuse
- Pour la robustesse et reproductibilité : Forêt Aléatoire**
 - Stabilité intermédiaire : écart-type **5,38%**, plus prévisible que GBM
 - Sensibilité très élevée : **97,50%** pour la détection d'AML
 - Performance raisonnable : **56,79%** d'accuracy
 - Limitation** : même biais que la régression logistique (spécificité = 0,0%)

Recommandation contextuelle : Dans un contexte médical où l'équilibre entre sensibilité et spécificité est crucial, le **Gradient Boosting** s'impose malgré sa variabilité, étant la seule méthode capable de détecter les deux classes (spécificité = 30,00%).

6.6 Conclusion de l'étude comparative

Cette analyse comparative démontre plusieurs principes fondamentaux de l'apprentissage automatique en grande dimension, quantifiés par nos résultats empiriques :

1. Nuancement de la supériorité des méthodes d'ensemble : Contrairement aux attentes théoriques, la régression logistique Ridge atteint la meilleure accuracy globale (**58,21%**), dépassant le gradient

boosting (**57,86%**) et la forêt aléatoire (**56,79%**). Cependant, cette supériorité masque un biais complet (spécificité = 0,0%) qui limite son applicabilité pratique.

2. Importance critique du choix méthodologique quantifiée : L'écart de performance modéré (**1,42 points** entre meilleure et moins bonne accuracy) souligne que le choix doit privilégier les caractéristiques qualitatives : - **Stabilité** : Régression logistique (écart-type = 2,26%) >> Forêt aléatoire (5,38%) >> Gradient boosting (20,55%) - **Équilibre classes** : Seul le gradient boosting atteint une spécificité non-nulle (**30%**) - **Sensibilité** : Forêt aléatoire (**97,50%**) \approx Régression logistique (**100%**) > Gradient boosting (**78%**)

3. Validation empirique des limites en grande dimension : Le ratio $p \gg n$ ($1000/72 \approx 14$) impose des contraintes sévères confirmées quantitativement : - Aucune méthode ne dépasse **60% d'accuracy**, illustrant la difficulté intrinsèque - Le déséquilibre des classes (58,3% AML / 41,7% ALL) amplifie les biais algorithmiques - La régularisation optimale ($\lambda = 0,18$ pour Ridge) privilégie la stabilité au détriment de la discrimination

4. Adaptation contextuelle nécessaire : Dans notre contexte médical spécifique (classification leucémies), le **gradient boosting** s'impose comme choix optimal malgré sa variabilité (**écart-type 20,55%**), étant la seule méthode capable de discrimination bi-classe, critère essentiel en diagnostic différentiel.

Cette étude confirme empiriquement qu'en contexte $p \gg n$, **aucune méthode n'est universellement optimale** : le choix doit intégrer les contraintes métier (ici, la nécessité de détecter les deux classes de leucémies) au-delà des seules métriques de performance globale. La validation croisée rigoureuse révèle des nuances cruciales masquées par les moyennes agrégées.