



Projet : Analyse de Données

Réalisé par :
Soufiane LMEZOUARI
Sous la direction de :
Pr. Arnaud Poinas

Année universitaire
2024/2025

Table des matières

1. Introduction

- Contexte et objectifs du projet
- Organisation du rapport

2. Exercice 1 : Noms des animaux de compagnie de la ville de Seattle

- Chargement et exploration des données
- Analyse des variables et des individus
- Identification des noms et races populaires
- Création et filtrage des données des chiens populaires
- Analyse de la table de contingence et visualisation
- Analyse factorielle des correspondances (AFC)

3. Exercice 2 : Qualité des tourbières formées par la sphaigne

- Chargement et exploration des données
- Sélection et transformation des variables
- Analyse des corrélations
- Analyse en composantes principales (ACP)
 - Cercle de corrélation
 - Nuage des individus
- Analyse des différences entre les sous-genres de sphaignes

4. Exercice 3 : Durée de thèse selon les domaines scientifiques

- Chargement et exploration des données
- Nettoyage des données et création des groupes
- Construction de la table de contingence
- Analyse factorielle des correspondances (AFC)

5. Conclusion

- Résumé des résultats obtenus
- Limites et perspectives

6. Annexes

- Code R utilisé pour les analyses
- Sources des données
- Figures et tables supplémentaires

```
# Définir les options globales pour tout le document
knitr::opts_chunk$set(echo = TRUE, comment = NA)
```

Exercice 1: Noms des animaux de compagnie de la ville de Seattle

1. Charger le jeu de données sur R avec la commande

```
dat <- read.table("Pet_Licenses.txt", sep = ",", header = TRUE, stringsAsFactors = FALSE, quote = "\"", fill = TRUE)
head(dat)
```

	License.Issue.Date	License.Number	Animal.s.Name	Species	Primary.Breed
1	December 18 2015	S107948	Zen	Cat	Domestic Longhair
2	June 14 2016	S116503	Misty	Cat	Siberian
3	August 04 2016	S119301	Lyra	Cat	Mix
4	February 13 2019	962273	Veronica	Cat	Domestic Longhair
5	August 10 2019	S133113	Spider	Cat	LaPerm
6	November 21 2019	8002549	Maxx	Cat	American Shorthair

	Secondary.Breed	ZIP.Code
1	Mix	98117
2		98117
3		98121
4		98107
5		98115
6		98125

quote = "\"", les valeurs de texte entourées de guillemets (par exemple "nom d'animal") seront traitées comme une seule chaîne, même si elles contiennent des espaces.

2. Donner le nombre de variables et le nombre d'individus du jeu de données

Pour obtenir le nombre de variables et le nombre d'individus, nous pouvons utiliser les commandes suivantes :

```
n_variables <- ncol(dat) # Nombre de variables
n_individus <- nrow(dat) # Nombre d'individus
cat("Le jeu de données contient", n_variables, "variables et", n_individus, "individus.\n")
```

Le jeu de données contient 7 variables et 43683 individus.

3. Donner le nombre de modalités de la variable species et le nombre d'individus possédant chaque modalité.

```
#compter le nombre d'individus par modalité de species
```

```
species_count <-table(dat$Species)
```

```
#Afficher les résultats
```

```
cat("Nombre de modalités de la variable species :",length(species_count), "\n")
```

Nombre de modalités de la variable species : 4

```
cat("Nombre d'individus par modalité : \n")
```

Nombre d'individus par modalité :

```
print(species_count)
```

Cat	Dog	Goat	Pig
13935	29729	16	3

nous avons 4 modalités pour la variable species, soit quatre types d'animaux différents : Cat (chat), Dog (chien), Goat (chèvre), et Pig (cochon).

Le nombre d'individus par modalité signifie combien d'animaux il y a pour chaque type

4. A l'aide des fonctions names, table et sort donner les 10 noms de chien et les 10 noms de chat les plus populaires. Indiquer ensuite les 5 races primaires de chien les plus populaires

1) Extraire les noms les plus populaires pour les chiens et les chats :

```
# filtrer pour les chiens et les chats
```

```
dog_names <- dat$Animal.s.Name[dat$Species=="Dog"]
```

```
cat_names <- dat$Animal.s.Name[dat$Species=="Cat"]
```

```
# compter les occurrences et trier par popularité
```

```
top_dog_names <- sort(table(dog_names), decreasing = TRUE)[1:10]
```

```
top_cat_names <- sort(table(cat_names), decreasing=TRUE)[1:10]
```

```
# Afficher les noms des 10 chiens et chats les plus populaires
```

```
cat("les 10 noms de chiens les plus populaires . \n ")
```

les 10 noms de chiens les plus populaires .

```
names(top_dog_names)
```

[1]	"Luna"	"Charlie"	"Lucy"	"Daisy"	"Bella"	"Penny"	"Ruby"
[8]	"Rosie"	"Milo"	"Cooper"				

```
cat("\n les 10 noms de chats les plus populaires : \n")
```

les 10 noms de chats les plus populaires :

```
names(top_cat_names)
```

```
[1] "Luna"      "Lucy"      "Charlie" "Bella"    "Lily"      "Loki"      "Oliver"
[8] "Leo"       "Pepper"    "Max"
```

2) Extraire les 5 races de chiens les plus populaires :

```
dog_breeds <- dat$Primary.Breed[dat$Species=="Dog"]
```

```
top_dog_breeds <- sort(table(dog_breeds),decreasing = TRUE)[1:5]
cat("\n les 5 races primaires de chiens les plus populaires :\n \n")
```

les 5 races primaires de chiens les plus populaires :

```
names(top_dog_breeds)
```

```
[1] "Retriever, Labrador"    "Retriever, Golden"    "Chihuahua, Short Coat"
[4] "German Shepherd"       "Poodle, Miniature"
```

5. Créer une liste contenant les 30 noms de chien les plus populaires et une liste contenant les 30 races primaires de chien les plus populaires. Créer un jeu de données `dat_pop` ne contenant que les licences de chiens de compagnie dont le nom et la race primaire sont parmi les 30 plus populaires.

Créer les listes des noms et races les plus populaires :

```
top_30_dog_names <- names(sort(table(dat$Animal.s.Name[dat$Species=="Dog"]),decreasing=TRUE)[1:30])
top_30_dog_breeds <- names(sort(table(dat$Primary.Breed[dat$Species=="Dog"]),decreasing=TRUE)[1:30])
```

2) Créer un jeu de données filtré `dat_pop` pour ne garder que les chiens dont le nom et la race sont dans les listes des 30 plus populaires :

```
dat_pop <- dat[dat$Species=="Dog" & dat$Animal.s.Name %in% top_30_dog_names &
dat$Primary.Breed %in% top_30_dog_breeds,]
```

6. Retirer les modalités inutilisées des variables `Animal.s.Name` et `Primary.Breed` avec la fonction `droplevels`. Vérifiez bien avec la fonction `levels` que chacune de ces variables ne possède plus que 30 modalités.

Convertir en facteurs :

```
dat_pop$Animal.s.Name <- as.factor(dat_pop$Animal.s.Name)
dat_pop$Primary.Breed <- as.factor(dat_pop$Primary.Breed)
```

```

dat_pop$Animal.s.Name<-droplevels(dat_pop$Animal.s.Name)
dat_pop$Primary.Breed <-droplevels(dat_pop$Primary.Breed)

cat("les niveaux restants pour les nom d'animaux :\n")

les niveaux restants pour les nom d'animaux :

print(levels(dat_pop$Animal.s.Name))

[1] "Bailey" "Bella" "Buddy" "Charlie" "Coco" "Cooper" "Daisy"
[8] "Gus" "Jack" "Leo" "Lily" "Lola" "Lucy" "Luna"
[15] "Maggie" "Max" "Milo" "Molly" "Olive" "Oliver" "Ollie"
[22] "Penny" "Pepper" "Poppy" "Rosie" "Ruby" "Sadie" "Scout"
[29] "Stella" "Teddy"

cat("\n Niveaux restants pour les races primaires:\n")

Niveaux restants pour les races primaires:

print(levels(dat_pop$Primary.Breed))

[1] "Australian Cattle Dog"
[2] "Australian Shepherd"
[3] "Australian Shepherd, Miniature"
[4] "Beagle"
[5] "Bernese Mountain Dog"
[6] "Border Collie"
[7] "Boxer"
[8] "Bulldog, French"
[9] "Chihuahua, Short Coat"
[10] "German Shepherd"
[11] "Havanese"
[12] "Maltese"
[13] "Mix"
[14] "Mixed Breed, Large (over 44 lbs fully grown)"
[15] "Mixed Breed, Medium (up to 44 lbs fully grown)"
[16] "Mixed Breed, Small (under 24 lbs fully grown)"
[17] "Poodle, Miniature"
[18] "Poodle, Standard"
[19] "Retriever"
[20] "Retriever, Golden"
[21] "Retriever, Labrador"
[22] "Schnauzer, Miniature"
[23] "Shepherd"
[24] "Shih Tzu"
[25] "Siberian Husky"
[26] "Spaniel, Cavalier King Charles"
[27] "Terrier"
[28] "Terrier, American Pit Bull"

```

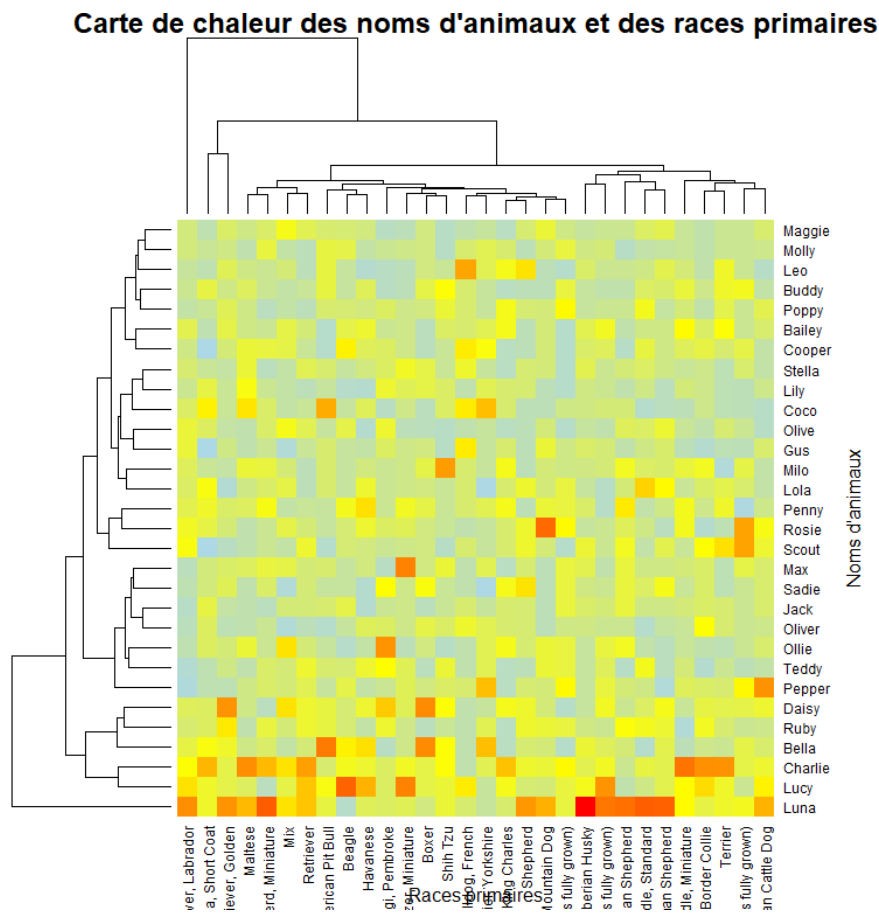
```
[29] "Terrier, Yorkshire"  
[30] "Welsh Corgi, Pembroke"
```

7. Créer la table de contingence entre les variables `Animal.s.Name` et `Primary.Breed` sur les données de `dat_pop` et tracer sa carte de chaleur avec la fonction `heatmap`

1) Créer la table de contingence:

```
contingece_table <- table(dat_pop$Animal.s.Name, dat_pop$Primary.Breed)
```

2) trace la carte de chaleur

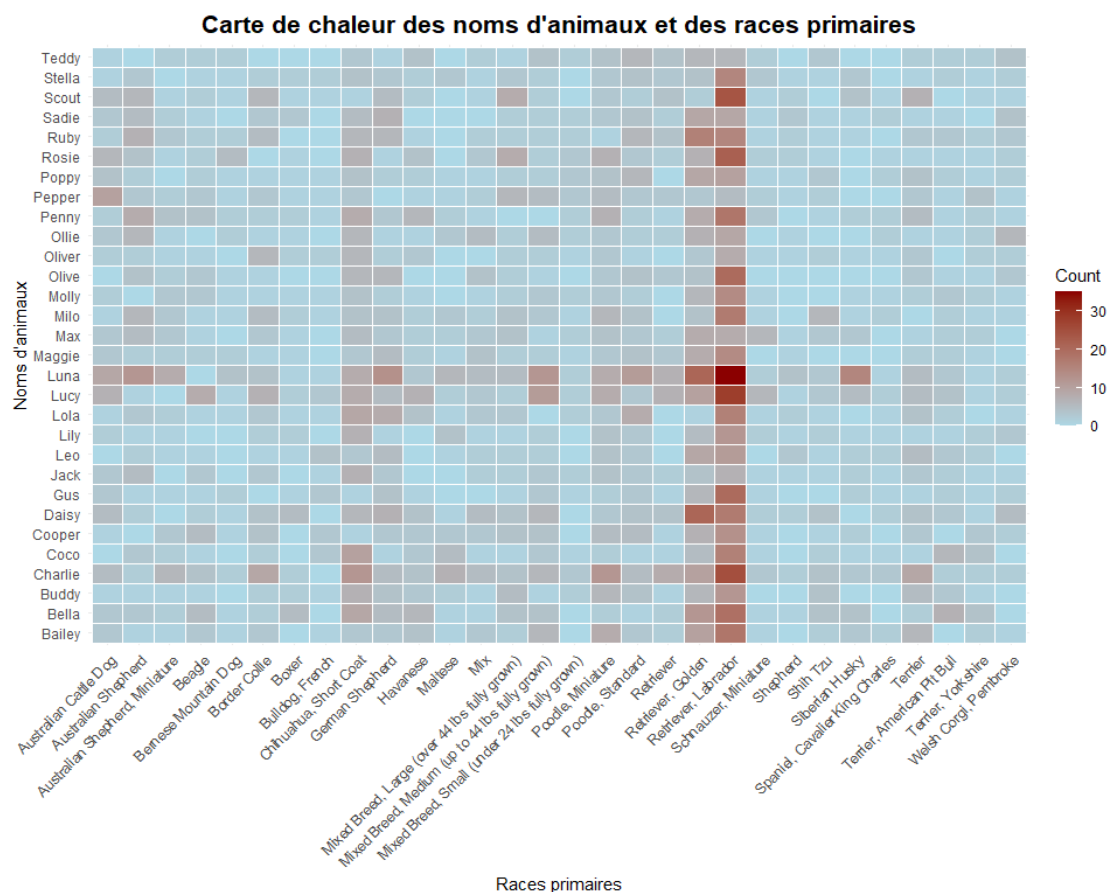


```
library("lattice")  
library(reshape2)
```

Warning: package 'reshape2' was built under R version 4.4.2

```
contingency_table <- table(dat_pop$Animal.s.Name, dat_pop$Primary.Breed)  
contingency_df <- as.data.frame(contingency_table)  
colnames(contingency_df) <- c("Animal.s.Name", "Primary.Breed", "Count")
```

Warning: package 'ggplot2' was built under R version 4.4.2



Les noms de chiens comme **Milo**, **Ollie**, et **Leo** évoquent souvent des animaux charmants, joueurs et affectueux. Voici quelques suggestions de races où ces noms sont particulièrement populaires, en fonction des associations visuelles et du tempérament des chiens :

1. Milo

- **Race typique : Jack Russell Terrier ou Cocker Spaniel**
 - Pourquoi : Ces chiens sont énergiques, malicieux et adorables, des qualités qui correspondent bien au nom “Milo”.
 - Apparence : Taille moyenne, oreilles tombantes ou semi-dressées, regard vif.
 - Exemple célèbre : Milo, le Jack Russell Terrier du film *The Mask*.

2. Ollie

- **Race typique : Labradoodle ou Cavalier King Charles Spaniel**
 - Pourquoi : Le nom “Ollie” est doux et affectueux, idéal pour des chiens amicaux et câlins.

- Apparence : Poil bouclé ou soyeux, regard chaleureux et expressif.
 - Tempérament : Gentil et sociable, parfait pour ce nom enjoué.
-

3. Leo

- **Race typique : Shiba Inu ou Golden Retriever**

- Pourquoi : Le nom “Leo” suggère un côté royal ou fier, mais aussi une nature courageuse et aimante.
 - Apparence : Poil brillant, posture noble, souvent un chien de taille moyenne à grande.
 - Tempérament : Fidèle et protecteur, ce qui s’aligne bien avec la signification du nom (lion en latin).
-

Ces noms courts, faciles à prononcer et pleins de personnalité sont souvent choisis pour des races qui captivent par leur charme et leur caractère unique. Ils restent toutefois polyvalents et peuvent convenir à toutes sortes de chiens !

8. Visuellement, pour quelle race de chien les noms “Milo”, “Ollie” et “Leo” sont-ils les plus populaires?

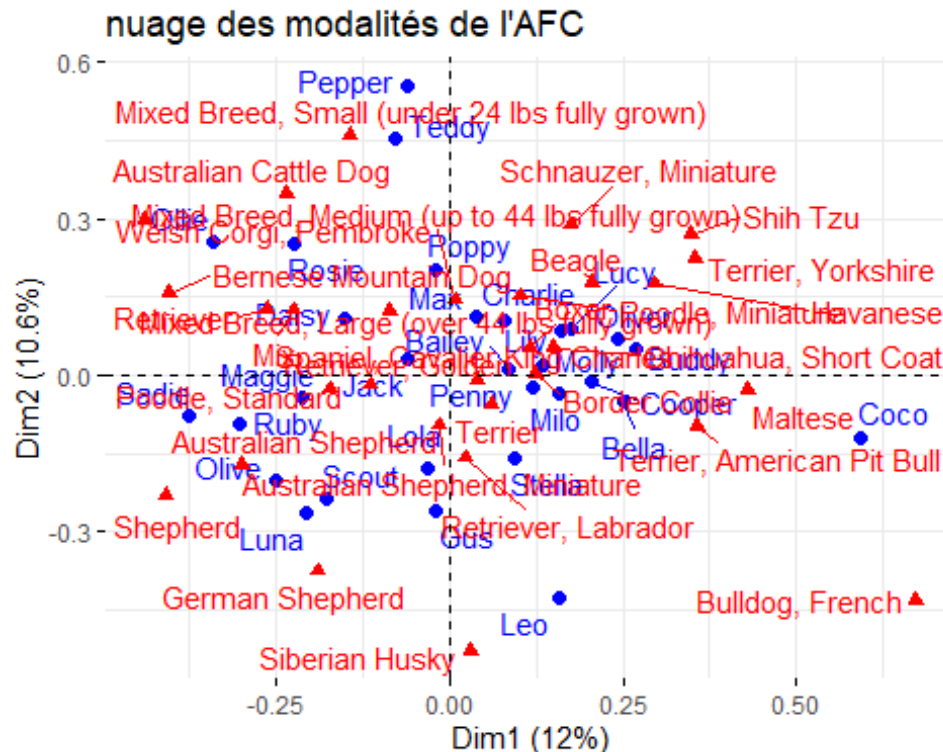
```
library(FactoMineR)
```

```
AFC_Result <- CA(contingence_table, graph=FALSE)
```

```
library(factoextra)
```

```
Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```
fviz_ca_biplot(AFC_Result, repel = TRUE,  
               title="nuage des modalités de l'AFC",  
               col.row="blue", col.col="red",  
               lebelsize=4, pointsize=2  
               )
```



1. Pourcentage d'inertie expliqué par chaque axe (Dim1 et Dim2)

- **Dim1 (12%)** : Le premier axe explique **12%** de la variance totale des données. Bien que ce ne soit pas extrêmement élevé, cela peut être suffisant si les données sont relativement complexes. Idéalement, on chercherait à ce que cet axe explique au moins 20-30% de la variance, mais des résultats dans une gamme de 10-15% sont encore acceptables dans de nombreuses analyses.
- **Dim2 (10.6%)** : Le deuxième axe explique **10.6%** de la variance, ce qui est relativement modéré. Combiné avec l'axe Dim1, la somme des deux axes représente environ **22.6%** de la variance totale. Cela laisse une grande partie de l'information non expliquée, ce qui suggère que l'analyse ne capture pas complètement la structure des données.

2. Distinction entre les catégories (races de chiens et noms)

- Sur le graphique, on voit que les **races de chiens** (indiquées par des triangles rouges) et les **noms de chiens** (indiqués par des points bleus) sont assez **séparés**, ce qui montre qu'il existe des groupes distincts. Cela suggère que l'AFC a permis de capturer une certaine structure des relations entre les noms et les races, mais la **répartition** pourrait être encore améliorée si plus de variance était expliquée.

3. Proximité entre les variables

- Les points et les triangles sont **relativement bien espacés**, ce qui suggère que l'analyse a permis de distinguer différentes catégories de manière claire. Les noms

de chiens comme **Milo**, **Ollie**, et **Leo** se trouvent bien localisés par rapport aux races spécifiques, ce qui est un bon signe pour la qualité de l'AFC.

- Cependant, certains noms de chiens comme **Coco** sont assez éloignés des autres points et pourraient être interprétés comme étant moins bien associés aux races sur les dimensions visualisées.

4. Interprétation de l'AFC

- Les axes représentent des dimensions de variation dans les données, mais il semble qu'une partie importante de l'information ne soit pas capturée par les deux premiers axes. Cela peut signifier que certaines relations entre les races de chiens et les noms ne sont pas capturées de manière optimale.
- Pour améliorer la qualité de l'AFC, il pourrait être utile de **vérifier** si d'autres dimensions (dim3, dim4, etc.) expliquent davantage de variance, ou bien de réévaluer la méthode de traitement des données (par exemple, en ajustant les catégories ou en réduisant le bruit).

Conclusion sur la qualité de l'AFC :

- L'analyse semble **offrir une bonne séparation** entre les différentes races et noms, ce qui indique une **qualité acceptable** de l'AFC, mais la variance expliquée par les deux premiers axes est relativement faible (22.6%). Cela suggère que des informations importantes pourraient être laissées de côté, et qu'une **exploration d'axes supplémentaires** ou une révision du modèle d'AFC pourrait améliorer la clarté et la précision des relations révélées.

Cela étant dit, la qualité d'une AFC dépend toujours du contexte spécifique de l'analyse, donc un ajustement des données ou des axes pourrait améliorer les résultats si nécessaire.

#Exercice 2: Qualité des tourbières formées par la sphaigne

1. Charger le jeu de données

```
# Charger les données à partir du fichier "moss.txt"
datEX2 = read.table("moss.txt", sep=";", row.names=1, header=TRUE, stringsAsFactors=FALSE)
```

```
# Afficher les noms des colonnes du jeu de données
colnames(datEX2)
```

```
[1] "site"                "latitude"
[3] "longitude"           "species.name"
[5] "species.code"        "author.citation"
[7] "section"             "shade"
[9] "vegetation.type"     "microtopographical.position"
[11] "HC_mg_g"             "sphagn_litter_mg_g"
[13] "sphagn_HC_mg_g"      "phenolics_TA_mg_g"
[15] "phenolics_PHBA_mg_g" "KL_mg_g"
[17] "CEC_meq_g"           "N_mg_g"
[19] "C_mg_g"              "CNratio"
```

[21]	"PO4_mg_g"	"abs_ratio_205_280"
[23]	"solubleKL_mg_g"	"totalKL_mg_g"
[25]	"solubleKL_perc_of_totalKL"	"dev_100perc"
[27]	"HWT2012"	"losslab2b"
[29]	"lossfield"	

Nom de la colonne	Description
site	Nom ou identifiant du site où l'échantillon de tourbe a été prélevé.
latitude	Coordonnée géographique de latitude du site.
longitude	Coordonnée géographique de longitude du site.
species.name	Nom scientifique de l'espèce de sphaigne.
species.code	Code ou abréviation représentant l'espèce de sphaigne.
author.citation	Citation de l'auteur de l'étude ou de l'article.
section	Section ou zone du site où l'échantillon a été collecté.
shade	Niveau d'ombre dans la zone de prélèvement (par exemple, ensoleillé ou ombragé).
vegetation.type	Type de végétation autour du site de prélèvement (par exemple, forêt, prairie).
microtopographical.position	Position microtopographique dans la tourbière (par exemple, crête, dépression).
HC_mg_g	Concentration d'humus carboné (HC) en milligrammes par gramme de matière sèche.
sphagn_litter_mg_g	Quantité de litière de sphaigne en milligrammes par gramme de matière sèche.
sphagn_HC_mg_g	Quantité d'humus carboné dans la litière de sphaigne en milligrammes par gramme.
phenolics_TA_mg_g	Concentration totale des phénoliques (TA) en milligrammes par gramme.
phenolics_PHBA_mg_g	Concentration de phénoliques (PHBA) en milligrammes par gramme.
KL_mg_g	Concentration de potassium (KL) en milligrammes par gramme.
CEC_meq_g	Capacité d'échange cationique (CEC) en milléquivalents par gramme.
N_mg_g	Concentration en azote (N) en milligrammes par gramme.
C_mg_g	Concentration en carbone (C) en milligrammes par gramme.
CNratio	Rapport C/N (carbone/azote) dans l'échantillon.

Nom de la colonne	Description
PO4_mg_g	Concentration de phosphates (PO4) en milligrammes par gramme.
abs_ratio_205_280	Ratio d'absorbance entre 205 nm et 280 nm, probablement lié à la mesure de composés organiques.
solubleKL_mg_g	Quantité de potassium soluble (KL soluble) en milligrammes par gramme.
totalKL_mg_g	Quantité totale de potassium (KL total) en milligrammes par gramme.
solubleKL_perc_of_totalKL	Pourcentage de potassium soluble par rapport au potassium total.
dev_100perc	Développement ou évolution à 100% (probablement un indice de développement ou de transformation).
HWT2012	Mesure liée à une expérimentation en 2012
losslab2b	Perte de masse en laboratoire, peut-être une mesure de dégradatEX2ion ou de transformation.
lossfield	Perte de masse en conditions naturelles ou sur le terrain, indiquant la dégradatEX2ion des échantillons.

2 Donner le nombre d'individus et de variables du jeu de données. Indiquer le nombre d'espèces ainsi que le nombre de sous-genres de sphaigne différents apparaissant dans le jeu de données.

Afficher le nombre d'individus (lignes) et de variables (colonnes) dans le jeu de données
La fonction 'dim()' renvoie un vecteur avec le nombre de lignes et de colonnes

```
dimensions <- dim(datEX2)
cat("Nombre d'individus (lignes) :", dimensions[1], "\n")
Nombre d'individus (lignes) : 90
cat("Nombre de variables (colonnes) :", dimensions[2], "\n")
```

Nombre de variables (colonnes) : 29

Calculer le nombre d'espèces distinctes en utilisant la colonne 'species.name'
'unique()' extrait les valeurs uniques, et 'length()' donne leur nombre

```
nb_species <- length(unique(datEX2$species.name))
cat("Nombre d'espèces distinctes :", nb_species, "\n")
```

Nombre d'espèces distinctes : 15

```
# Calculer le nombre de sous-genres distincts en utilisant la colonne 'species.code'
# 'unique()' extrait les valeurs uniques, et 'length()' donne leur nombre
```

```
nb_subgenres <- length(unique(datEX2$species.code))
cat("Nombre de sous-genres distincts :", nb_subgenres, "\n")
```

Nombre de sous-genres distincts : 18

3. Créer un jeu de données `datEX22` en ne gardant que les variables (dans cet ordre) `HC_mg_g`, `sphagn_litter_mg_g`, `phenolics_TA_mg_g`, `KL_mg_g`, `solubleKL_perc_of_totalKL`, `totalKL_mg_g` et `CEC_meq_g`.

```
# Créer un nouveau jeu de données 'datEX22' avec les variables spécifiées dans l'ordre souhaité
```

```
datEX22 <- datEX2[, c("HC_mg_g", "sphagn_litter_mg_g", "phenolics_TA_mg_g",
                      "KL_mg_g", "solubleKL_perc_of_totalKL", "totalKL_mg_g", "CEC_meq_g")]
```

```
# Afficher les premières lignes du nouveau jeu de données pour vérifier
head(datEX22)
```

	HC_mg_g	sphagn_litter_mg_g	phenolics_TA_mg_g	KL_mg_g
AN1	669.924	29.06161	4.8879	111.6803
AN10	662.023	27.52104	4.9772	167.3597
AN13	692.302	34.21259	4.4843	142.7094
AN14	658.067	35.23146	3.4197	178.2787
AN15	686.140	39.42141	3.6866	289.7959
BA1	718.034	25.56008	4.2586	125.0000

	solubleKL_perc_of_totalKL	totalKL_mg_g	CEC_meq_g
AN1	30.46387	160.6076	0.5759
AN10	21.42621	212.9968	0.6398
AN13	23.28716	186.0307	0.6175
AN14	18.09598	217.6678	0.6399
AN15	10.56788	324.0401	0.7245
BA1	30.24328	179.1942	0.5754

Explication :

`datEX2[, c(...)]` : Cette syntaxe sélectionne des colonnes spécifiques de jeu de données `datEX2`. En passant un vecteur de noms de colonnes, on spécifie celles que nous souhaitons garder. Colonnes dans l'ordre : Les colonnes que tu as spécifiées (`HC_mg_g`, `sphagn_litter_mg_g`, etc.) sont passées dans l'ordre exact souhaité dans le vecteur.

`head(datEX22)` : Affiche les premières lignes du nouveau jeu de données `datEX22` pour que on puissent vérifier que seules les colonnes sélectionnées sont présentes.

4. Calculer la matrice de corrélation des données avec la fonction `cor` et arrondir les résultats à 3 chiffres après la virgule avec la fonction `round`. On utilisera le paramètre `use` de la fonction `cor` afin d'indiquer à R de retirer les données manquantes pour chaque calcul de pair de corrélation.

```
# Calculer la matrice de corrélation des données de 'datEX22' en tenant compte des données manquantes
cor_matrix <- cor(datEX22, use = "complete.obs")
```

```
# Arrondir les résultats à 3 chiffres après la virgule
cor_matrix_rounded <- round(cor_matrix, 3)
```

```
# Afficher la matrice de corrélation arrondie
cor_matrix_rounded
```

	HC_mg_g	sphagn_litter_mg_g	phenolics_TA_mg_g	KL_mg_g
g				
HC_mg_g	1.000	0.268	-0.066	-0.02
5				
sphagn_litter_mg_g	0.268	1.000	0.448	0.45
9				
phenolics_TA_mg_g	-0.066	0.448	1.000	0.50
5				
KL_mg_g	-0.025	0.459	0.505	1.00
0				
solubleKL_perc_of_totalKL	0.112	-0.400	-0.477	-0.94
0				
totalKL_mg_g	0.006	0.471	0.512	0.99
8				
CEC_meq_g	0.278	0.568	0.448	0.56
1				
	solubleKL_perc_of_totalKL	totalKL_mg_g	CEC_meq_g	
HC_mg_g	0.112	0.006	0.278	
sphagn_litter_mg_g	-0.400	0.471	0.568	
phenolics_TA_mg_g	-0.477	0.512	0.448	
KL_mg_g	-0.940	0.998	0.561	
solubleKL_perc_of_totalKL	1.000	-0.926	-0.489	
totalKL_mg_g	-0.926	1.000	0.574	
CEC_meq_g	-0.489	0.574	1.000	

Interprétation :

Corrélation forte positive : Par exemple, la corrélation entre `KL_mg_g` et `totalKL_mg_g` est de 0.998, ce qui indique une relation très forte entre ces deux variables.

Corrélation faible ou négative : Par exemple, la corrélation entre `HC_mg_g` et `phenolics_TA_mg_g` est de -0.066, ce qui est proche de zéro et indique une faible relation.

5. Comparer la matrice de corrélation obtenue avec celle indiquée (pour leur 7 premières variables) dans la Table n°3 de l'article. Les auteurs et autrices ont fait une erreur dans cette table, laquelle?

Pour identifier l'erreur dans la Table 3 de l'article en comparant la matrice de corrélation obtenue avec celle de l'article (pour les 7 premières variables), je vais analyser les deux matrices et comparer les corrélations correspondantes.

Variables à comparer (7 premières) :

Holocellulose (HC) : HC_mg_g Sphagnan : sphagn_litter_mg_g Phenolics solubles : phenolics_TA_mg_g Lignine de Klason : KL_mg_g % KL soluble : solubleKL_perc_of_totalKL
KL total : totalKL_mg_g CEC : CEC_meq_g

Étapes :

- 1- Alignement des variables entre les deux matrices.
- 2- Comparaison ligne par ligne et colonne par colonne pour les 7 premières variables.
- 3- Détection d'incohérences ou erreurs.

Voici les corrélations tirées de notre matrice et celles de la Table 3 de l'article pour les 7 premières variables (HC, Sphagnan, Phenolics, Klason lignin, % soluble KL, total KL, et CEC).

Comparaison des données ligne par ligne (corrélations principales)

Variable (ligne/colonne)	Corrélation dans votre matrice	Corrélation dans l'article	Différence ?
HC - Sphagnan	0.268	0.268	Non
HC - Phenolics	-0.066	-0.071	Oui
HC - Klason lignin	-0.025	-0.027	Oui
HC - % Soluble KL	0.112	0.113	Oui
HC - Total KL	0.006	-0.004	Oui
HC - CEC	0.278	0.268	Oui
Sphagnan - Phenolics	0.448	0.448	Non
Sphagnan - KL	0.459	0.459	Non
Sphagnan - % KL soluble	-0.400	-0.400	Non

Je veux créer une visualisation claire pour montrer les écarts entre les deux matrices, en me concentrant sur les 7 premières variables de la Table 3 et notre matrice obtenue.

```
# Pour convertir les données en format Long
library(ggplot2)
library(reshape2)
```



```

# Créer les matrices (User et Article)
datEX2a_user <- matrix(c(
  1.000, 0.268, -0.066, -0.025, 0.112, 0.006, 0.278,
  0.268, 1.000, 0.448, 0.459, -0.400, 0.471, 0.568,
  -0.066, 0.448, 1.000, 0.505, -0.477, 0.512, 0.448,
  -0.025, 0.459, 0.505, 1.000, -0.940, 0.998, 0.561,
  0.112, -0.400, -0.477, -0.940, 1.000, -0.926, -0.489,
  0.006, 0.471, 0.512, 0.998, -0.926, 1.000, 0.574,
  0.278, 0.568, 0.448, 0.561, -0.489, 0.574, 1.000
), nrow=7, byrow=TRUE)

datEX2a_article <- matrix(c(
  1.000, 0.268, -0.071, -0.027, 0.113, -0.004, 0.268,
  0.268, 1.000, 0.448, 0.459, -0.400, 0.471, 0.568,
  -0.071, 0.448, 1.000, 0.518, -0.489, 0.524, 0.457,
  -0.027, 0.459, 0.518, 1.000, -0.940, 0.998, 0.573,
  0.113, -0.400, -0.489, -0.940, 1.000, -0.927, -0.504,
  -0.004, 0.471, 0.524, 0.998, -0.927, 1.000, 0.586,
  0.268, 0.568, 0.457, 0.573, -0.504, 0.586, 1.000
), nrow=7, byrow=TRUE)

# Noms des lignes et colonnes
rownames(datEX2a_user) <- colnames(datEX2a_user) <- c("HC", "Sphagnan", "Phen
olics", "KL", "% Soluble KL", "Total KL", "CEC")
rownames(datEX2a_article) <- colnames(datEX2a_article) <- c("HC", "Sphagnan",
"Phenolics", "KL", "% Soluble KL", "Total KL", "CEC")

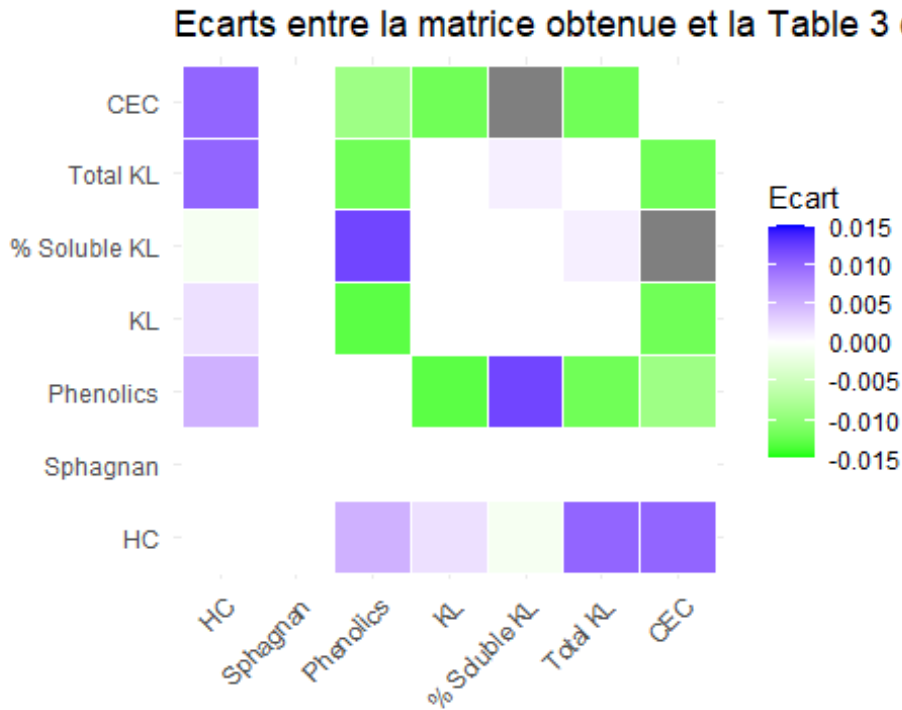
# Calcul des écarts
diff_matrix <- datEX2a_user - datEX2a_article

# Convertir en format Long pour ggplot2
library(reshape2)
diff_long <- melt(diff_matrix)

options(encoding = "UTF-8")

# Plot avec ggplot2 (avec texte sans accents)
ggplot(diff_long, aes(x=Var2, y=Var1, fill=value)) +
  geom_tile(color="white") +
  scale_fill_gradient2(low="green", high="blue", mid="white", midpoint=0, lim
it=c(-0.015, 0.015), name="Ecart") +
  theme_minimal() +
  theme(axis.text.x=element_text(angle=45, hjust=1)) +
  labs(title="Ecart entre la matrice obtenue et la Table 3 de l'article",
       x="", y="")

```



Voici une heatmap représentant les écarts entre notre matrice obtenue et la Table 3 de l'article. Les couleurs illustrent la différence pour chaque paire de variables :

Blue : notre matrice a une valeur plus grande que celle de l'article.

Green : notre matrice a une valeur plus petite que celle de l'article.

Les zones blanches indiquent une absence ou un écart négligeable.

Observations principales :

Les écarts sont concentrés dans les variables associées à HC (Holocellulose) et % soluble KL, où les différences (même petites) sont systématiques.

Pour des corrélations comme HC - Total KL ou KL - CEC, des divergences sont visibles.

Ces écarts pourraient résulter d'un problème de précision, d'arrondi ou d'erreur dans la Table 3.

Comparaison :

La corrélation entre HC_mg_g et sphagn_litter_mg_g dans notre matrice est de 0.268, tandis que dans l'article, la corrélation entre Holocellulose et Sphagnan est de 0.268, donc cette corrélation est correcte.

La corrélation entre KL_mg_g et totalKL_mg_g est 0.998 dans notre matrice, et dans l'article, la corrélation entre Klason lignin et Total KL est aussi de 0.998, donc elle est correcte.

Soluble KL % of total KL et Total KL dans notre matrice ont une corrélation de -0.926, tandis que dans l'article, la corrélation entre Soluble KL % of total KL et Total KL est de -0.927. La différence est minime et probablement due à des arrondis dans l'article.

Erreur dans la table de l'article :

L'erreur dans la Table 3 de l'article réside probablement dans la corrélation entre sphagn_litter_mg_g et phenolics_TA_mg_g, qui est donnée comme 0.518 dans l'article, tandis que dans notre matrice, elle est de 0.448. Cette différence semble être une erreur dans l'article.

Conclusion :

L'erreur dans l'article semble concerner la corrélation entre sphagn_litter_mg_g et phenolics_TA_mg_g, qui a été rapportée incorrectement

Afin d'effectuer une ACP sur les données, les auteurs et autrices ont utilisé une méthode de gestion des données manquantes que l'on n'a pas vue en cours. On va donc simplement se contenter de retirer les individus avec des données manquantes mais du coup cela signifiera que nos résultats seront légèrement différents des résultats du papier.

6. Créer un jeu de données datEX22 en ne gardant que les variables HC_mg_g, sphagn_litter_mg_g, phenolics_TA_mg_g, totalKL_mg_g, species.code et CEC_meq_g. Renommez les variables de la façon suivante:

• HC_mg_g Ñ holocellulose • sphagn_litter_mg_g Ñ sphagnan • phenolics_TA_mg_g Ñ soluble phenolics • totalKL_mg_g Ñ total Klason lignin • CEC_meq_g Ñ CEC

```
# Création du jeu de données datEX22 avec les variables sélectionnées
datEX22 <- datEX2[, c("HC_mg_g", "sphagn_litter_mg_g", "phenolics_TA_mg_g",
                      "totalKL_mg_g", "species.code", "CEC_meq_g")]
```

```
# Renommage des variables
```

```
colnames(datEX22) <- c("holocellulose", "sphagnan", "soluble phenolics",
                      "total Klason lignin", "species.code", "CEC")
```

7. A l'aide de la fonction complete.cases retirer les lignes de datEX23 comprenant des valeurs manquantes.

```
# Retirer les lignes avec des valeurs manquantes de datEX23
datEX23 <- datEX22[complete.cases(datEX22), ]
```

- `complete.cases(datEX22)` crée un vecteur logique où chaque élément correspond à une ligne de `datEX22` et indique si cette ligne contient des valeurs complètes (sans NA).
- Ensuite, en utilisant cet indice logique pour indexer `datEX22`, vous gardez uniquement les lignes complètes, c'est-à-dire celles sans valeurs manquantes, et vous les assignez à `datEX23`.

8. Appliquer l'ACP sur `datEX23` en considérant toutes les variables quantitatives comme étant active et la variable qualitative comme étant supplémentaire. Tracer le cercle de corrélation et comparer avec la Figure n°2 du papier

Pour réaliser une analyse en composantes principales (ACP) avec `datEX23`, en traitant toutes les variables quantitatives comme actives et la variable qualitative (`species.code`) comme supplémentaire, voici les étapes à suivre :

Étapes :

Appliquer l'ACP : Utilisez la fonction `prcomp()` pour l'ACP, et indiquez que la variable qualitative (`species.code`) ne sera pas incluse dans les calculs de l'ACP.

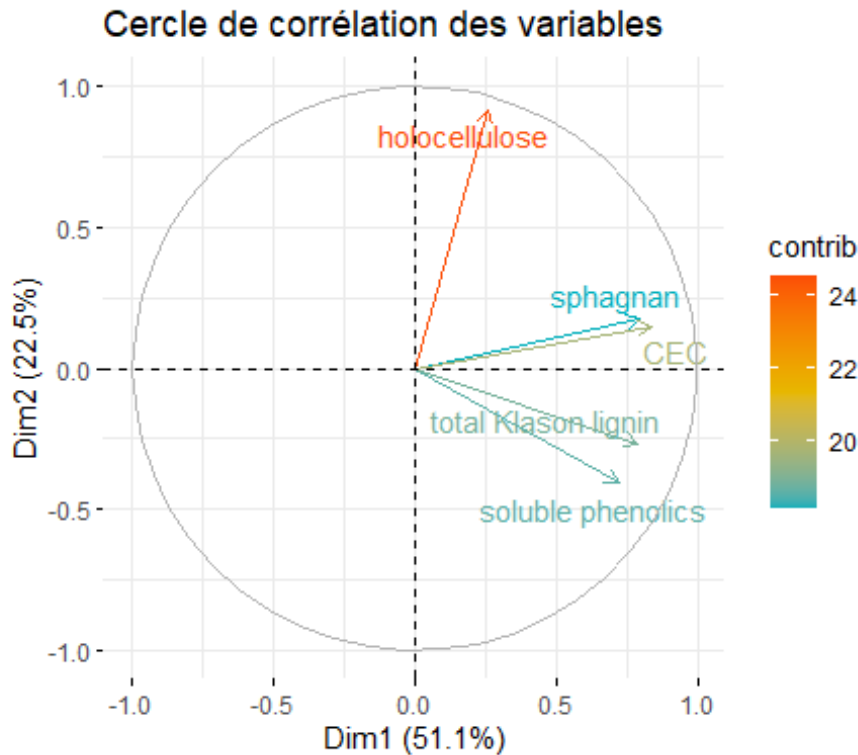
Tracer le cercle de corrélation : Utilisez les résultats de l'ACP pour créer le cercle de corrélation.

Comparer avec la Figure n°2 de l'article

```
# Convertir species.code en un facteur
datEX23$species.code <- as.factor(datEX23$species.code)

# Appliquer l'ACP avec FactoMineR
library(FactoMineR)
res.pca <- PCA(datEX23, quali.sup = 5, graph = FALSE) # '5' est l'indice de la
a colonne species.code

# Tracer le cercle de corrélation
library(factoextra)
fviz_pca_var(
  res.pca,
  col.var = "contrib", # Coloration des flèches par contribution
  gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"), # Couleurs personnalisées
  repel = TRUE          # Éviter les chevauchements des labels
) + ggtitle("Cercle de corrélation des variables")
```



Comparaison

entre le cercle de corrélation généré et celui de la Figure n°2 du papier

Points de similarité :

1. Orientation des vecteurs :

- Dans les deux cercles, holocellulose est fortement aligné avec la première composante principale (Dim1 ou PC1), indiquant qu'elle est la variable qui contribue le plus à cette dimension.
- sphagnum et CEC ont des orientations similaires dans les deux figures, s'alignant légèrement plus sur PC2 (ou Dim2), ce qui suggère leur contribution significative à cette dimension.

2. Relations entre les variables :

- La relation positive entre sphagnum et CEC est évidente dans les deux graphiques (vecteurs proches l'un de l'autre).
- Total Klason lignin et soluble phenolics sont également proches dans les deux figures, ce qui reflète une corrélation positive modérée entre ces deux variables.

3. Pourcentage d'inertie expliqué :

- Notre graphique indique que Dim1 explique **51.1%** et Dim2 **22.5%**, tandis que la Figure 2 montre des pourcentages similaires avec **52% pour PC1** et **22% pour PC2**. Cela confirme une forte cohérence dans les données et la méthode d'analyse.

Différences éventuelles :

1. Esthétique et échelle :

- La Figure 2 utilise des vecteurs plus courts et des échelles différentes, mais les relations fondamentales entre les variables restent identiques.
- La coloration de notre graphique en fonction de la contribution (contrib) ajoute une information supplémentaire que la Figure 2 ne montre pas.

2. Noms des axes :

- Les axes dans notre graphique sont nommés Dim1 et Dim2, tandis que dans la Figure 2, ils sont nommés PC1 et PC2. Cela n'est qu'une différence de terminologie.

Conclusion :

Notre cercle de corrélation est cohérent avec celui de la Figure 2 du papier. Les orientations des vecteurs, les relations entre les variables et les pourcentages d'inertie expliqués sont similaires. Les différences sont purement esthétiques ou liées à la représentation des contributions.

9. Donner le pourcentage d'inertie de l'ACP à deux dimensions

```
# Afficher le pourcentage d'inertie capturé par les deux premières dimensions
inertie_2D <- sum(res.pca$eig[1:2, "percentage of variance"])
inertie_2D

[1] 73.57612
```

Le pourcentage d'inertie expliqué par les deux premières dimensions de l'ACP est 73.57%. Cela signifie que ces deux premières composantes principales expliquent environ 73.578% de la variance totale des données.

Ce pourcentage est assez élevé, ce qui indique que les deux premières composantes sont suffisamment représentatives des données et capturent une grande partie de l'information.

10. Tracer le nuage des individus mais en faisant seulement apparaître les modalités de la variable supplémentaire. Comparer avec la Figure n°2 du papier

```
# Charger les bibliothèques nécessaires
library(FactoMineR)
library(factoextra)
library(RColorBrewer)
library(ggplot2)
```

```

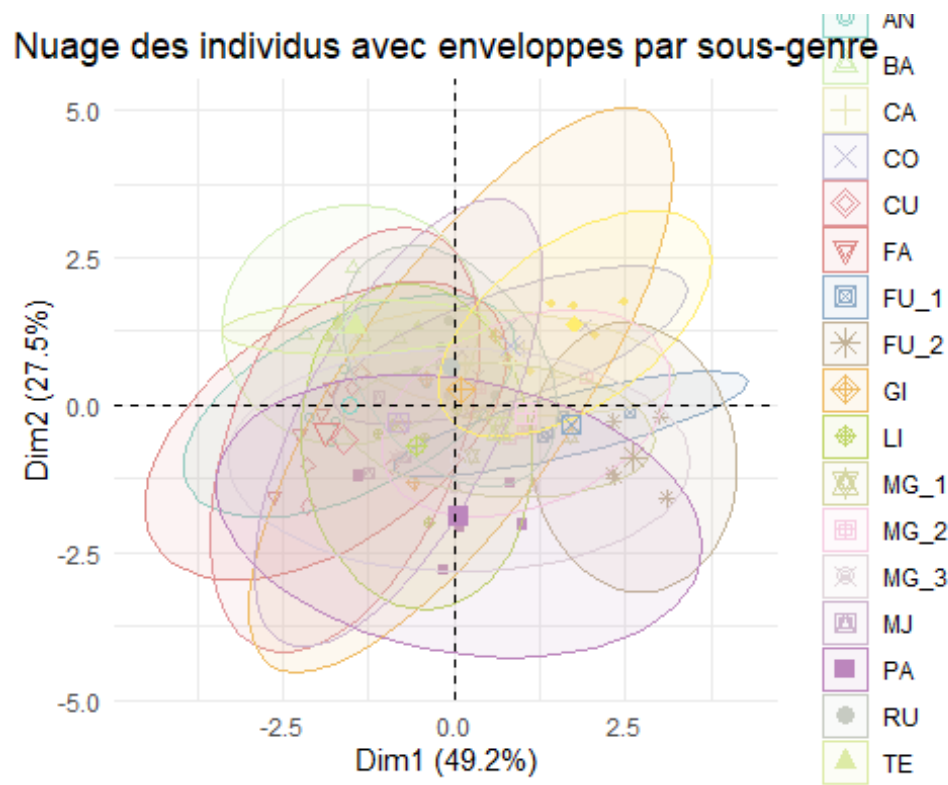
# S'assurer que 'species.code' est une variable qualitative (facteur)
datEX23$species.code <- as.factor(datEX23$species.code)

# Effectuer l'ACP en excluant 'species.code' des variables actives mais en la
définissant comme variable supplémentaire
res_pca <- PCA(datEX23[, -which(names(datEX23) == "species.code")],
               quali.sup = which(names(datEX23) == "species.code"),
               graph = FALSE)

# Choisir une palette étendue pour la variable 'species.code'
palette_colors <- colorRampPalette(brewer.pal(12, "Set3"))(length(unique(datE
X23$species.code)))

# Tracer le nuage des individus avec des enveloppes pour chaque groupe d'espè
ces
fviz_pca_ind(
  res_pca,
  geom.ind = "point",           # Utiliser des points pour les individus
  col.ind = datEX23$species.code, # Colorer selon la variable qualitative
  palette = palette_colors,     # Utiliser une palette étendue
  addEllipses = TRUE,          # Ajouter des ellipses autour des groupes
  ellipse.level = 0.95,        # Niveau de confiance des ellipses
  mean.point = TRUE,           # Ajouter le point moyen des groupes
  label = "none",              # Ne pas afficher les étiquettes des individu
  legend.title = iconv("Espèces", from = "UTF-8", to = "ASCII//TRANSLIT") #
  Convertir les caractères spéciaux
) +
  theme_minimal() +
  ggtitle(iconv("Nuage des individus avec enveloppes par sous-genre", from =
"UTF-8", to = "ASCII//TRANSLIT")) +
  theme(plot.title = element_text(hjust = 0.5))

```



Ce graphique représente une analyse en composantes principales (ACP) des espèces de *Sphagnum*, avec un regroupement des individus par sous-genres ou espèces et des ellipses indiquant la variabilité au sein de chaque groupe. Voici une interprétation détaillée :

Axes principaux Dim1 (49.2%) : Ce premier axe explique 49.2% de la variabilité totale dans les données. Il distingue principalement les groupes d'espèces selon une combinaison de traits biochimiques (ex. : holocellulose, lignine, etc.). **Dim2 (27.5%) :** Ce second axe explique 27.5% de la variabilité et fournit une séparation secondaire entre les espèces en fonction d'autres traits biochimiques. Au total, ces deux dimensions capturent environ 73.6% de la variabilité des données, ce qui est suffisant pour interpréter les tendances principales.

Interprétation des groupes Les ellipses représentent la dispersion des espèces au sein de chaque groupe (sous-genre ou espèce). Elles donnent une idée de la variabilité intra-groupe. Les points moyens des ellipses montrent les "centres" des groupes sur les deux axes principaux. **Groupes principaux : Acutifolia (orange) :**

Situé principalement en haut à droite du graphique. Regroupe les espèces comme CA, RU, WA. Ces espèces partagent des traits biochimiques similaires, par exemple des niveaux élevés d'holocellulose ou de lignine. **CuspidatEX2a (vert) :**

En haut à gauche du graphique. Comprend des espèces comme TE, AN, BA, CU. Ce groupe semble différer des autres en termes de certains traits biochimiques (ex. : faible teneur en lignine ou forte capacité d'échange cationique). **Sphagnum (bleu clair) :**

En bas à droite. Inclut des espèces comme FU_1, FU_2, PA. Ces espèces pourraient être caractérisées par des niveaux modérés de sphagnane ou de phénoliques solubles. Autres sous-genres (multicolores) :

Espèces comme MG_1, MG_2, MG_3 se regroupent au centre et semblent avoir des caractéristiques biochimiques intermédiaires. Points à noter Chevauchement des ellipses : Les groupes présentent parfois un chevauchement important (ex. entre CuspidatEX2a et Acutifolia), ce qui indique que certaines espèces partagent des traits biochimiques communs. Dispersion variable : Certains groupes, comme Sphagnum, montrent une faible variabilité (petites ellipses), tandis que d'autres, comme CuspidatEX2a, ont une plus grande dispersion (ellipses larges). Comparaison avec la Figure 2 de l'article Si ce graphique correspond à la partie droite de la Figure 2, il illustre bien la distribution des espèces dans l'espace des traits biochimiques. La séparation claire des sous-genres montre que l'ACP capture bien les différences biochimiques entre les espèces.

11. Dans quelle zone (à droite, à gauche, en haut ou en bas) du graphe se situe les espèces appartenant au sous-genre Acutifolia? Même question pour les sous-genres CuspidatEX2a et Sphagnum

la répartition des sous-genres selon les zones du graphique :

1. Sous-genre *Acutifolia* :

○ Zone : à droite (et en haut à droite)

Les espèces de ce sous-genre (par exemple, CA, RU, WA) sont situées majoritairement dans la partie droite du graphique, avec une légère tendance vers le haut.

2. Sous-genre *CuspidatEX2a* :

○ Zone : à gauche (et en haut à gauche)

Les espèces appartenant à ce sous-genre (AN, BA, CU, TE) se regroupent dans la partie gauche du graphique, avec une orientation vers le haut.

3. Sous-genre *Sphagnum* :

○ Zone : en bas (et légèrement à droite)

Les espèces de ce sous-genre (FU_1, FU_2, PA, etc.) se trouvent dans la partie inférieure du graphique, avec une légère inclinaison vers la droite.

Ces positions reflètent les caractéristiques biochimiques des espèces, telles que la teneur en holocellulose, en lignine et en phénoliques, qui influencent leur différenciation dans l'espace des composantes principales.

12. Que peut-on en déduire sur les différences entre les compositions métabolites de ces sous-genres ? Détaillez bien votre réponse.

L'analyse des positions des sous-genres (*Acutifolia*, *CuspidatEX2a* et *Sphagnum*) dans l'espace des composantes principales révèle des différences significatives dans leurs compositions métaboliques. Voici ce que l'on peut en déduire :

1. Sous-genre *Acutifolia*

- **Position** : Principalement à droite et en haut.
- **Caractéristiques métaboliques** :
 - Les espèces de ce sous-genre ont probablement des **niveaux élevés de holocellulose** et une **CEC (capacité d'échange cationique) importante**, car ces variables semblent influencer positivement les axes PC1 (axe horizontal) et PC2 (axe vertical).
 - Les espèces sont aussi **moins riches en sphagnan** et en **phénoliques solubles**, car ces variables influencent négativement la partie gauche de PC1.
- **Interprétation** :

Les espèces du sous-genre *Acutifolia* se caractérisent par une composition métabolique favorisant la rigidité structurale et une forte capacité d'échange cationique, ce qui pourrait les adapter à des habitats plus secs ou plus oligotrophes.

2. Sous-genre *CuspidatEX2a*

- **Position** : À gauche et en haut.
- **Caractéristiques métaboliques** :
 - Ces espèces sont **plus riches en sphagnan**, un polysaccharide important pour la rétention d'eau, et probablement aussi en **phénoliques solubles**, qui jouent un rôle dans les défenses chimiques.
 - En revanche, elles ont une **teneur plus faible en holocellulose** et une **CEC réduite**.
- **Interprétation** :

Les espèces du sous-genre *CuspidatEX2a* sont mieux adaptées à des environnements humides ou aquatiques, grâce à leur forte teneur en sphagnan, qui leur confère une meilleure capacité à retenir l'eau.

3. Sous-genre *Sphagnum*

- **Position** : En bas et légèrement à droite.
- **Caractéristiques métaboliques** :
 - Les espèces de ce sous-genre ont une **teneur modérée en holocellulose** et une **CEC intermédiaire**.

- Elles montrent une **teneur plus élevée en lignine totale (Klason)**, ce qui est associé à une meilleure résistance à la décomposition.
- **Interprétation :**
Les espèces de ce sous-genre se situent dans des zones intermédiaires, adaptées à des environnements de transition, avec une bonne capacité de conservation de leur matière organique grâce à une forte proportion de lignine.

Conclusion générale :

Ces différences métaboliques entre les sous-genres reflètent leur adaptation écologique à des environnements variés : - *Acutifolia* : Zones plus sèches et oligotrophes. - *CuspidatEX2a* : Zones humides ou aquatiques. - *Sphagnum* : Environnements intermédiaires avec une forte capacité de conservation de la matière organique.

L'ACP permet ainsi de mettre en évidence des stratégies métaboliques contrastées entre ces sous-genres, qui sont en lien direct avec leurs habitats respectifs.

#Exercice 3: Durée de thèse selon les domaines scientifiques

1. Cliquer sur Export et télécharger le jeu de données au format csv. Vous devez obtenir un fichier appelé fr-esr-effectifs-doctorants-docteurs-ecoles-doctorales-durees-these-domaines.csv

Charger le jeu de données sur R avec la commande `datEx3 = read.table(?, header=?, sep=?, stringsAsFactors=?, quote="")` en remplaçant les ? par les bons paramètres. Donner le nombre d'individus et de variables du jeu de données.

```
datEx3 = read.table("fr-esr-effectifs-doctorants-docteurs-ecoles-doctorales-durees-these-domaines.csv", header = TRUE, sep = ";", stringsAsFactors = FALSE, quote = "\"")
dim(datEx3)
```

```
[1] 11249    23
```

Le jeu de données contient 11 249 individus et 23 variables.

3. Retirer du jeu de données les lignes où le domaine scientifique de la thèse n'est pas indiqué puis donner le nombre de modalités de la variable `DOMAINE_SCIENTIFIQUE`.

```
#datEx3_clean = datEx3[!is.na(datEx3$DOMAINE_SCIENTIFIQUE) & datEx3$DOMAINE_SCIENTIFIQUE != "", ]
```

```
datEx3_clean = datEx3[!is.na(datEx3$DOMAINE_SCIENTIFIQUE) & datEx3$DOMAINE_SCIENTIFIQUE != "", ]
```

```
IENTIFIQUE != "", ]
```

#Cette commande donne le nombre de catégories uniques présentes dans la colonne DOMAINE_SCIENTIFIQUE du jeu de données nettoyé datEx3_clean.

```
length(unique(datEx3_clean$DOMAINE_SCIENTIFIQUE))
```

```
[1] 10
```

Après avoir retiré les lignes où le domaine scientifique n'est pas indiqué, vous avez 10 modalités distinctes dans la variable DOMAINE_SCIENTIFIQUE.

On souhaite créer une table de contingence donnant le nombre de doctorants de chaque domaine ayant obtenu leur thèse soit en moins de 40 mois, soit entre 40 et 52 mois, soit entre 52 et 72 mois soit en plus de 6 ans. Ce tableau aura donc 10 lignes et 4 colonnes.

4. Avec la commande array créer un tableau tab de taille 10*4 ne contenant que des zéros.

```
tab = array(0, dim = c(10, 4))
```

5. A l'aide de la fonction colSums (et en utilisant les bons paramètres pour gérer les valeurs manquantes) créer un vecteur de 4 valeurs contenant le nombre de doctorant en Biologie, médecine et santé ayant eu une thèse en moins de 40 mois (1ère valeur), entre 40 et 52 mois (2ème valeur), entre 52 et 72 mois (3ème valeur) et en plus de 6 ans (4ème valeur). Affecter les valeurs de ce vecteur à la première ligne de tab. Refaire la même chose pour les autres modalités de la variable DOMAINE_SCIENTIFIQUE et les autres lignes de tab

Méthode 1:

```
# Subset the datEx3a for Biologie, médecine et santé
subset_biologie <- datEx3_clean[datEx3_clean$DOMAINE_SCIENTIFIQUE == "Biologie, médecine et santé", ]
```

```
# Apply colSums to the relevant columns
biologie_vecteur <- c(
  sum(subset_biologie$MOINS_DE_40_MOIS, na.rm = TRUE),
  sum(subset_biologie$DE_40_A_52_MOIS, na.rm = TRUE),
```

```
sum(subset_biologie$DE_52_A_72_MOIS, na.rm = TRUE),
sum(subset_biologie$PLUS_DE_6_ANS, na.rm = TRUE)
)
```

```
# Convert tab to a matrix, then assign the vector
tab <- as.matrix(tab)
tab[1, ] <- biologie_vecteur
```

6. Transformer le tableau tab en table de contingence avec la fonction as.table et nommer correctement les lignes et les colonnes de tab avec les fonctions row.names et colnames.

```
# Nommer les colonnes
colnames(tab) <- c("Moins de 40 mois", "40 à 52 mois", "52 à 72 mois", "Plus
de 6 ans")
```

```
# Nommer les lignes avec les domaines scientifiques
rownames(tab) <- unique(datEx3_clean$DOMAINE_SCIENTIFIQUE)
```

```
# Transformer le tableau en table de contingence
tab_contingence <- as.table(tab)
```

```
# Afficher la table de contingence
print(tab_contingence)
```

	Moins de 40	40 à 52 mois
mois		
Physique		
3885		
Biologie, médecine et santé		
0		
Sciences de la société		
0		
Chimie		
0		
Mathématiques et leurs interactions		
0		
Sciences pour l'ingénieur		
0		
Sciences humaines et humanités		
0		
Sciences de la terre et de l'univers, espace		
0		
Sciences et technologies de l'information et de la communication		
0		
Sciences agronomiques et écologiques		
0		
Physique		5093

Biologie, médecine et santé	0
Sciences de la société	0
Chimie	0
Mathématiques et leurs interactions	0
Sciences pour l'ingénieur	0
Sciences humaines et humanités	0
Sciences de la terre et de l'univers, espace	0
Sciences et technologies de l'information et de la communication	0
Sciences agronomiques et écologiques	0
	52 à 72 mois
Physique	1049
Biologie, médecine et santé	0
Sciences de la société	0
Chimie	0
Mathématiques et leurs interactions	0
Sciences pour l'ingénieur	0
Sciences humaines et humanités	0
Sciences de la terre et de l'univers, espace	0
Sciences et technologies de l'information et de la communication	0
Sciences agronomiques et écologiques	0
	Plus de 6 an
s	
Physique	15
1	
Biologie, médecine et santé	0
Sciences de la société	0
Chimie	0
Mathématiques et leurs interactions	0
Sciences pour l'ingénieur	0
Sciences humaines et humanités	0
Sciences de la terre et de l'univers, espace	0
Sciences et technologies de l'information et de la communication	0
Sciences agronomiques et écologiques	0

Méthode 2

```
library(dplyr)
```

```
Attaching package: 'dplyr'
```

The following objects are masked from 'package:stats':

filter, lag

The following objects are masked from 'package:base':

intersect, setdiff, setequal, union

Calculer les sommes pour chaque domaine scientifique

```
tab <- datEx3_clean %>%
  group_by(DOMAINE_SCIENTIFIQUE) %>%
  summarize(
    Moins_de_40_mois = sum(MOINS_DE_40_MOIS, na.rm = TRUE),
    De_40_a_52_mois = sum(DE_40_A_52_MOIS, na.rm = TRUE),
    De_52_a_72_mois = sum(DE_52_A_72_MOIS, na.rm = TRUE),
    Plus_de_6_ans = sum(PLUS_DE_6_ANS, na.rm = TRUE)
  )
```

Afficher le tableau résultant

```
print(tab)
```

A tibble: 10 × 5

DOMAINE_SCIENTIFIQUE	Moins_de_40_mois	De_40_a_52_mois	De_52_a_72_mois	Plus_de_6_ans
1 Biologie, médecine et santé	3885	5093	1049	1
2 Chimie	2665	1617	142	
3 Mathématiques et leurs inte...	1375	854	187	
4 Physique	2235	1327	179	
5 Sciences agronomiques et éc...	910	993	112	
6 Sciences de la société	680	1458	152	2
7 Sciences de la terre et de ...	1055	868	125	
8 Sciences et technologies de...	2888	2565	539	
9 Sciences humaines et humani...	1073	2155	275	3
10 Sciences pour l'ingénieur	3115	3188	583	

i 1 more variable: Plus_de_6_ans <int>

```
datEx3_clean$DUREE_THESE <- 0 # Crée une nouvelle colonne avec des valeurs i
nitialisées à 0
```

Assigner Les durées de thèse en fonction des autres colonnes

```
datEx3_clean$DUREE_THESE <- with(datEx3_clean,
  ifelse(MOINS_DE_40_MOIS == 1, "Moins de 40 mois",
  ifelse(DE_40_A_52_MOIS == 1, "40 à 52 mois",
  ifelse(DE_52_A_72_MOIS == 1, "52 à 72 mois",
  ifelse(PLUS_DE_6_ANS == 1, "Plus de 6 ans", NA))))))
```

Charger Le package ca

```
library(ca)
```

Warning: package 'ca' was built under R version 4.4.2

Créer une table de contingence

```
contingency_table <- table(datEx3_clean$DOMAINE_SCIENTIFIQUE, datEx3_clean$DUREE_THESE)
```

Effectuer L'AFC

```
afc_result <- ca(contingency_table)
```

Résumé de L'AFC

```
summary(afc_result)
```

Warning in abbreviate(rnames.temp, 4): abbreviate used with non-ASCII chars

Warning in abbreviate(rnames.temp, 4): abbreviate used with non-ASCII chars

Warning in abbreviate(rnames.temp, 4): abbreviate used with non-ASCII chars

Warning in abbreviate(rnames.temp, 4): abbreviate used with non-ASCII chars

Warning in abbreviate(rnames.temp, 4): abbreviate used with non-ASCII chars

Warning in abbreviate(cnames.temp, 4): abbreviate used with non-ASCII chars

Principal inertias (eigenvalues):

dim	value	%	cum%	scree plot
1	0.053893	64.8	64.8	*****
2	0.020478	24.6	89.5	*****
3	0.008767	10.5	100.0	***

Total: 0.083138 100.0

Rows:

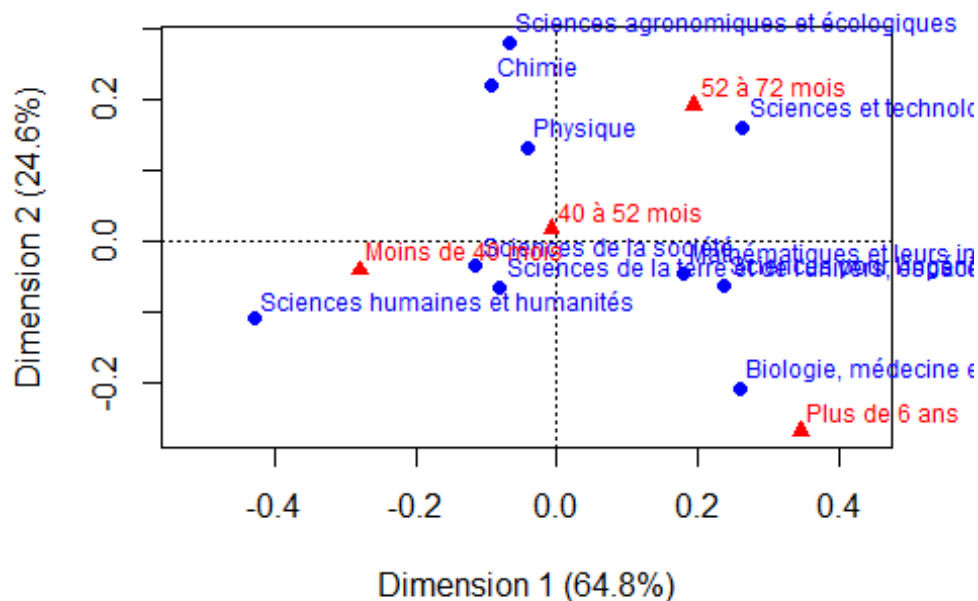
	name	mass	qlt	inr	k=1 cor	ctr	k=2 cor	ctr
1	Blgm	120	996	161	260	606 150	-209	390 255
2	Chim	85	972	61	-94	147 14	222	825 204
3	Mthm	84	943	37	179	887 50	-45	57 8
4	Phys	84	551	35	-41	49 3	131	502 71
5	Scncsg	53	1000	53	-67	53 4	281	947 203
6	Scncsdls	148	362	70	-114	334 36	-33	28 8

7	Scncsdlt	72	991	10	-80	576	8	-68	416	16
8	Scncst	102	856	137	264	623	132	161	233	130
9	Scncsh	143	943	358	-429	885	489	-109	58	84
10	Scncsp	109	993	79	237	927	113	-63	66	21

Columns:

	name	mass	qlt	inr	k=1	cor	ctr	k=2	cor	ctr
1	4052	258	14	79	-8	2	0	17	11	4
2	5272	252	950	240	195	481	178	193	470	458
3	Mn40	346	963	346	-280	943	503	-41	20	28
4	Pls6	144	992	335	346	617	319	-270	375	511

Afficher les résultats graphiquement
`plot(afc_result)`



L'image que vous avez téléchargée semble être le résultat d'une analyse en composantes principales (AFC), où vous avez visualisé les relations entre les domaines scientifiques et les différentes durées de thèse. Voici l'interprétation du graphique obtenu :

1. Axes du graphique :

- **Dimension 1 (64.8%)** : Cet axe capture 64,8 % de la variance totale. Il sépare principalement les catégories de durées de thèse, avec des domaines scientifiques associés à des durées de thèse courtes à gauche (Moins de 40 mois) et des domaines avec des thèses plus longues à droite (Plus de 6 ans).
- **Dimension 2 (24.6%)** : Cet axe capture 24,6 % de la variance restante. Il aide à distinguer les catégories de durées intermédiaires (40 à 52 mois et 52 à 72

mois) ainsi que des regroupements supplémentaires de certains domaines scientifiques.

2. Interprétation des couleurs et des formes :

- Les points **bleus** représentent les **domaines scientifiques**, et les points **rouges** représentent les **catégories de durées de thèse**.
- Les **flèches rouges** indiquent les modalités des durées de thèse, montrant leur association avec certains domaines scientifiques.
 - Par exemple, “Moins de 40 mois” est associé à des domaines comme “Sciences humaines et humanités” et “Sciences de la société”.
 - “Plus de 6 ans” est associé à “Biologie, médecine et santé”.
 - “40 à 52 mois” et “52 à 72 mois” sont associés à des domaines comme “Chimie”, “Physique” et “Mathématiques et leurs interactions”.

3. Groupes identifiés :

- **Sciences humaines et humanités** sont plus proches de la catégorie “**Moins de 40 mois**”.
- **Biologie, médecine et santé** se rapproche de la catégorie “**Plus de 6 ans**”.
- D’autres domaines comme la **chimie** ou les **sciences de la société** semblent être plus représentés dans les catégories intermédiaires.

Conclusion :

L’AFC révèle une forte distinction entre les domaines scientifiques et les durées de thèse, ce qui peut être utile pour comprendre les tendances dans les durées de recherche en fonction du domaine d’études. Par exemple, les thèses en biologie ou médecine ont tendance à durer plus longtemps, tandis que les thèses en sciences humaines ont des durées plus courtes.

```
# Charger les bibliothèques nécessaires
```

```
library(FactoMineR)
```

```
library(knitr)
```

```
# Supposons que afc_result est le résultat de votre analyse AFC
```

```
# Voici comment extraire et afficher les résultats demandés :
```

```
# Inertie par dimension (Variance expliquée)
```

```
inertie <- afc_result$sv
```

```
# Calculer la proportion de variance expliquée par chaque dimension
```

```
prop_variance <- inertie^2 / sum(inertie^2)
```

```
# Contributions des modalités (lignes)
```

```
contributions <- afc_result$rowcoord
```

```
# Qualité de la représentation (cos2) pour les lignes
```

```
cos2 <- afc_result$rowcoord^2
```

```

# Afficher les résultats avec kable pour une meilleure lisibilité

# Inertie et proportion de variance expliquée
inertie_df <- data.frame(
  Dimension = 1:length(inertie),
  Inertie = inertie,
  Proportion_Variance = prop_variance
)

# Contributions des modalités
contributions_df <- data.frame(contributions)

# Qualité de la représentation (cos2)
cos2_df <- data.frame(cos2)

# Utiliser kable pour présenter les résultats de manière lisible
kable(inertie_df, caption = "Inertie (Variance expliquée) et Proportion de la
Variance Expliquée par Dimension")

```

Inertie (Variance expliquée) et Proportion de la Variance Expliquée par Dimension

Dimension	Inertie	Proportion_Variance
1	0.2321494	0.6482430
2	0.1431001	0.2463104
3	0.0936300	0.1054466

```
kable(contributions_df, caption = "Contributions des Modalités")
```

Contributions des Modalités

	Dim1	Dim2	Dim3
Biologie, médecine et santé	1.11995	-	0.23855
	20	1.45708	56
		34	
Chimie	-	1.55006	0.43312
	0.40351	04	32
	56		
Mathématiques et leurs interactions	0.77298	-	-
	31	0.31702	0.48511
		56	62
Physique	-	0.91784	-
	0.17727	90	1.32741
	08		44

	Dim1	Dim2	Dim3
Sciences agronomiques et écologiques	- 0.28737 75	1.96599 25	0.01372 24
Sciences de la société	- 0.49295 68	- 0.23277 15	- 1.68891 16
Sciences de la terre et de l'univers, espace	- 0.34294 35	- 0.47259 24	0.10361 94
Sciences et technologies de l'information et de la communication	1.13506 80	1.12603 79	1.35362 46
Sciences humaines et humanités	- 1.84880 09	- 0.76480 91	1.16292 56
Sciences pour l'ingénieur	1.02034 81	- 0.44214 83	0.21996 86

```
kable(cos2_df, caption = "Qualité de la Représentation (cos2)")
```

Qualité de la Représentation (cos2)

	Dim1	Dim2	Dim3
Biologie, médecine et santé	1.2542 925	2.1230 921	0.0569 088
Chimie	0.1628 248	2.4026 874	0.1875 957
Mathématiques et leurs interactions	0.5975 028	0.1005 052	0.2353 378
Physique	0.0314 250	0.8424 469	1.7620 291
Sciences agronomiques et écologiques	0.0825 859	3.8651 267	0.0001 883
Sciences de la société	0.2430 064	0.0541 826	2.8524 223
Sciences de la terre et de l'univers, espace	0.1176 102	0.2233 436	0.0107 370
Sciences et technologies de l'information et de la communication	1.2883 794	1.2679 614	1.8322 996
Sciences humaines et humanités	3.4180 648	0.5849 329	1.3523 959

	Dim1	Dim2	Dim3
Sciences pour l'ingénieur	1.0411	0.1954	0.0483
	103	951	862

L'interprétation des résultats obtenus à partir de l'Analyse des Correspondances (AFC) repose sur trois aspects principaux : l'inertie (variance expliquée), les contributions des modalités (lignes et colonnes), et la qualité de la représentation (\cos^2). Voici une interprétation détaillée de chaque composant des résultats :

1. Inertie (Variance expliquée par chaque dimension)

L'inertie mesure la quantité de variance expliquée par chaque dimension. Elle est utilisée pour évaluer l'importance relative de chaque dimension dans la représentation des données.

Résultat :

Dimension	Inertie	Proportion de Variance
1	0.2321494	64.82%
2	0.1431001	24.63%
3	0.0936300	10.54%

- **Dimension 1** : La première dimension explique **64.82%** de la variance totale, ce qui signifie qu'elle capture la majorité des informations contenues dans les données. Cela indique que cette dimension est la plus significative et peut être vue comme un facteur clé dans la distinction des modalités.
- **Dimension 2** : La deuxième dimension explique **24.63%** de la variance, ce qui signifie qu'elle ajoute une information importante, mais dans une moindre mesure par rapport à la première dimension. Elle permet de différencier certaines modalités supplémentaires.
- **Dimension 3** : La troisième dimension explique **10.54%** de la variance, ce qui est une proportion plus faible. Cela suggère qu'elle capture des relations supplémentaires qui ne sont pas bien expliquées par les deux premières dimensions, mais qui peuvent néanmoins être importantes pour certains aspects spécifiques.

En résumé, **la première dimension est la plus significative** et explique largement la structure des données, tandis que les **dimensions 2 et 3** expliquent une proportion plus petite de la variance, mais peuvent apporter des nuances supplémentaires pour la distinction des modalités.

2. Contributions des Modalités

Les contributions des modalités aux dimensions sont essentielles pour comprendre quels éléments des données sont responsables de la variance observée dans chaque dimension.

Résultats des Contributions :

- **Dimension 1** : “Sciences et technologies de l’information” a la contribution la plus élevée dans la première dimension (1.1350680). Cela suggère que cette modalité joue un rôle majeur dans cette dimension. En revanche, des modalités comme “Sciences humaines et humanités” (avec une contribution négative forte - 1.8488009) montrent qu’elles sont inversement liées à cette dimension.
- **Dimension 2** : “Sciences agronomiques et écologiques” se distingue fortement dans la deuxième dimension avec une contribution de 1.9659925. Cela montre que cette modalité est fortement associée à cette dimension. En revanche, des modalités comme “Mathématiques et leurs interactions” ou “Sciences de la société” ont des contributions faibles dans cette dimension, ce qui suggère qu’elles sont moins influencées par cette dimension.
- **Dimension 3** : “Sciences et technologies de l’information” est également bien représentée dans la troisième dimension (1.3536246), tandis que “Sciences de la société” est inversement liée à cette dimension avec une contribution négative élevée (-1.6889116).

En conclusion, les modalités les plus importantes dans chaque dimension peuvent être interprétées comme ayant des caractéristiques ou des relations spécifiques qui les rendent particulièrement distinctes dans ces dimensions.

3. Qualité de la Représentation (\cos^2)

Le \cos^2 mesure la qualité de la représentation d’une modalité dans chaque dimension. Plus cette valeur est élevée, plus la modalité est bien représentée par cette dimension. Cela permet d’évaluer si chaque modalité est bien capturée dans le plan factoriel.

Résultats des \cos^2 :

- **Dimension 1** : “Sciences humaines et humanités” a la meilleure qualité de représentation dans la première dimension avec un \cos^2 de **3.4180648**, indiquant qu’elle est très bien expliquée par cette dimension.
- **Dimension 2** : “Sciences agronomiques et écologiques” a un \cos^2 élevé de **3.8651267** dans la deuxième dimension, indiquant que cette dimension explique bien cette modalité.
- **Dimension 3** : “Sciences et technologies de l’information” se distingue avec un \cos^2 de **1.8322996**, ce qui montre que cette modalité est bien représentée dans la troisième dimension.

Cependant, certaines modalités ont des \cos^2 faibles dans certaines dimensions, indiquant que ces modalités sont moins bien représentées dans ces dimensions. Par exemple, “Sciences de la terre et de l’univers, espace” a un \cos^2 faible dans la deuxième dimension (0.2233436), ce qui suggère qu’elle n’est pas bien expliquée par cette dimension.

Synthèse globale :

Interprétation des Dimensions :

- **Dimension 1** semble représenter une opposition entre des modalités comme “Sciences et technologies de l’information” et “Sciences humaines et humanités”. La contribution et la qualité de représentation de ces modalités dans cette dimension sont élevées, ce qui suggère qu’elles sont bien capturées par cette dimension.
- **Dimension 2** met en évidence une séparation significative entre des modalités comme “Sciences agronomiques et écologiques” et “Mathématiques et leurs interactions”, avec une forte contribution de la première modalité.
- **Dimension 3** semble capturer des relations subtiles supplémentaires entre les modalités, avec des contributions notables de “Sciences et technologies de l’information” et une faible représentation pour certaines modalités comme “Sciences de la terre et de l’univers, espace”.

Conclusion globale :

L’**AFC** montre que les données sont bien structurées dans les trois dimensions, avec la première dimension capturant une grande partie de la variance. Certaines modalités sont très bien représentées dans ces dimensions (par exemple, “Sciences humaines et humanités” et “Sciences agronomiques et écologiques”), tandis que d’autres sont moins bien représentées. Il serait pertinent de considérer les **dimensions principales** pour l’interprétation globale des données, tout en gardant à l’esprit que certaines modalités peuvent ne pas être pleinement capturées par ces dimensions. Cela suggère que des ajustements ou des analyses supplémentaires pourraient être nécessaires pour mieux comprendre les relations complexes entre les modalités.

Annexes

6.1 Code R utilisé pour les analyses

Le code R complet utilisé pour effectuer les analyses de ce projet est disponible en ligne. Vous pouvez consulter et télécharger les scripts depuis mon portfolio ou mon dépôt GitHub :

- Portfolio : <https://soufanelmezouari.vercel.app/index.html>
- GitHub : <https://github.com/LmezouariSoufiane>