

# Correspondance entre classes étiquettées et Clusters prédits par algorithme hongrois

Willy Kinfoussia, Aurélien Henriques

2023-2024 M2DS

## Introduction

Un problème notable dans l'apprentissage supervisé est la labélisation d'un grand nombre de données afin d'entraîner un modèle. En effet, Il arrive souvent que le nombre de données labélisées soit dérisoire comparé aux données dont l'on dispose au total. Certains algorithmes, par exemple l'algorithme SWAV, permettent de s'affranchir de ce problème en créant une représentation de classe par alteration faible des données labélisés, mais ils nécessitent un apprentissage contrastif afin de pouvoir identifier une classe.

Les algorithmes de clusterisation non-supervisé peuvent alors permettre de résoudre ce problème, en clusterisant les données pourvu que nous puissions identifier correctement les clusters aux étiquettes de classe. Dans ce cas, il suffirait alors de s'assurer de la bonne représentation du modèle des différentes classes via les clusters, et d'apposer aux données des clusters les étiquettes de classe issue de l'association cluster/classe.

Dans ce projet, nous proposons d'appliquer l'algorithme hongrois à ce problème. L'algorithme hongrois est un programme qui permet de trouver la meilleure association entre  $N$  agents et  $N$  tâches à l'aide de poids associés aux arrêtes entre agents et tâches. En pratique, il existe  $N^2$  poids d'arrête et  $N!$  façons différentes d'associer les agents aux tâches. La recherche de l'optimum des associations peut donc être pénible et utiliser l'algorithme hongrois permet de résoudre en partie ce problème, grâce une complexité polynômiale en  $N$ . Il repose sur une recherche du poids d'arrête minimal dans une matrice de coût qui contiendrait tous les poids des arrêtes entre agents et tâches (dans notre cas, entre cluster et classe). Nous utilisons ici cet algorithme ainsi qu'un algorithme naïf parcourant les  $N!$  possibilités

## Construction de la matrice de coût

Nous collectons les données MNIST, ainsi que les classes associées pour obtenir une matrice de coût. Pour ceci, nous effectuons une PCA à 50 dimensions sur les données puis nous clusterisons par k-nn les données obtenues en 10 clusters (le nombre de classe dans MNIST). Nous calculons alors les centres d'inertie des classes et des clusters puis l'on utilise la distance euclidienne entre chaque couple classe/cluster.

$$d(c, k) = \sqrt{\sum_{l=1}^{50} (c_l - k_l)^2}$$

Où : -  $c_l$  représente la valeur du centre d'inertie de la classe  $c$  en la  $l$ -ème dimension, -  $k_l$  représente la valeur du centre d'inertie du cluster  $k$  en la  $l$ -ème dimension,

```
mnist_train <- read.csv('mnist_test.csv')  
x_train <- mnist_train[, -1]
```

```

y_train <- mnist_train[, 1]

x_train_vec <- x_train / 255

constant_columns <- apply(x_train_vec, 2, function(col) length(unique(col)) == 1)
x_train_vec <- x_train_vec[, !constant_columns]

pca <- prcomp(x_train_vec, center = TRUE, scale. = TRUE, rank = 50)

x_train_pca <- predict(pca, x_train_vec)

```

```

library(stats)

k <- 10

kmeans_model <- kmeans(x_train_pca, centers = k)

clusters_centroid <- kmeans_model$centers

assignement_cluster <- kmeans_model$cluster

classes_centroid <- list()

```

## Optimisation

La matrice de coût est alors :

$$M = (d_{ck})$$

Où : - c et k varient entre 1 et 10.

La matrice  $X$  d'association classe/cluster de taille (10,10) est définie comme :

$$X_{ck} = \begin{cases} 1 & \text{si la classe } c \text{ et le cluster } k \text{ sont associés} \\ 0 & \text{sinon} \end{cases}$$

avec la contrainte qu'il ne doit y avoir qu'un seul 1 par ligne et des 0 sur le reste, correspondant à ne pas associer la classe c aux autres clusters.

$$\sum_{k=1}^{10} X_{ck} = 1$$

```

for (class_label in 0:(k - 1)) {

  class_data <- x_train_pca[y_train == class_label, ]

  centroid <- colMeans(class_data)

  classes_centroid[[as.character(class_label)]] <- centroid
}

```

```

cost_matrix_mnist <- matrix(0, nrow = k, ncol = k)

for (i in 0:k-1) {
  for (j in 1:k) {

    distance <- sqrt(sum((classes_centroid[[as.character(i)]] - clusters_centroid[j, ])^2))

    cost_matrix_mnist[i+1, j] <- distance
  }
}

cost_matrix_mnist <- lapply(1:k, function(i) as.vector(cost_matrix_mnist[i, ]))

#Génération de matrice de coût aléatoire

matrice_couts <- function(n) {

  cost_matrix_random <- matrix(rnorm(n * n, mean = 100, sd = 15), nrow = n, ncol = n)
  liste_matrice <- lapply(1:n, function(i) as.vector(cost_matrix_random[i, ]))

  return(liste_matrice)
}

```

Comme on souhaite aussi qu'un cluster corresponde à une ligne ( éventuellement il est possible de faire concorder plusieurs clusters à une classe et inversement en enlevant les contraintes sur la bijection de notre association), une seconde contrainte est qu'il ne doit y avoir qu'un 1 par ligne et des 0 sur le reste.

$$\sum_{c=1}^{10} X_{ck} = 1$$

La fonction de coût à minimiser est la somme des produits  $X_{kl} \times M_{kl}$  sur  $l$  et  $k$  :

$$\text{Coût}(\text{association } X) = \sum_{c=1}^{10} \sum_{k=1}^{10} X_{ck} \times M_{ck}$$

Cette fonction correspond au coût de l'association classe/cluster choisie. En la minimisant, on s'assure que l'association permet au maximum de réduire la distance séparant le cluster de la classe.

## Algorithme naïf

L'algorithme naïf permet de résoudre ce problème brutalement en calculant toutes les coûts possibles. C'est-à-dire calculer pour chaque bijection classes-cluster la somme des arrêtes associant un cluster à une classe et renvoyer la bijection qui minimise le coût. La complexité théorique de cet algorithme est  $N!$  où  $N$  est le nombre de classes/clusters. En effet, pour la classe 1, il y a 10 clusters possibles. Une fois un cluster choisi, on obtient le premier coefficient de coût. Il reste donc pour la classe 2 9 clusters possibles, 8 pour la classe 3, etc. Ceci donne bien  $N!$  associations possibles donc  $N!$  somme des coûts possibles, ce qui donne la complexité de cet algorithme.

## Explication détaillée de l'algorithme :

```
library(Rcpp)
```

```
## Warning: le package 'Rcpp' a été compilé avec la version R 4.3.3
```

```
sourceCpp("NaiveAlgorithme.cpp")
source("NaiveAlgorithme.R")
```

### Fonction NaiveAlgorithme :

Cette fonction résout le problème d'affectation en générant toutes les permutations possibles de l'assignation et en trouvant celle avec le coût minimal. Elle prend en entrée la matrice d'association. Elle initialise le vecteur minAssignment pour stocker l'assignation optimale avec le coût minimal et le coût minimal lui-même à une valeur maximale. Ensuite, elle initialise un vecteur currentAssignment pour stocker l'assignation courante. Elle génère tous les indices de permutation de 0 à n-1. Puis, elle itère à travers toutes les permutations possibles : Elle construit l'assignation courante en utilisant les indices de permutation. Elle calcule le coût de l'assignation courante en utilisant la fonction calculateCost. Elle met à jour l'assignation optimale et le coût minimal si le coût de l'assignation courante est inférieur au coût minimal actuel. Après avoir exploré toutes les permutations, elle renvoie une paire contenant l'assignation avec le coût minimal et le coût minimal lui-même.

### Fonction calculateCost :

Cette fonction calcule le coût total d'une affectation donnée en fonction de la matrice d'association et de l'assignation. Elle initialise le coût total à zéro. Ensuite, elle parcourt chaque ligne de la matrice d'association. Pour chaque ligne, elle ajoute le coût de l'élément correspondant dans la matrice d'association selon l'assignation donnée. Enfin, elle retourne le coût total calculé.

Exemple d'utilisation des packages :

- Package Rcpp :

```
cost_matrix <- list(
  c(1, 1, 5),
  c(5, 9, 7),
  c(8, 1, 9)
)
naive_adj_matrix <- NaiveAlgorithme(cost_matrix)
naive_adj_matrix
```

```
## [[1]]
## [1] 1 0 0
##
## [[2]]
## [1] 0 0 1
##
## [[3]]
## [1] 0 1 0
```

- Package R :

```
# Exemple d'utilisation :
# Créer une matrice de coût (exemple)
costMatrix <- matrix(c(1, 1, 5, 5, 9, 7, 8, 1, 9), nrow = 3)

# Appeler l'algorithme naïf
NaiveAlgorithme_R(costMatrix)
```

```
##      [,1] [,2] [,3]
## [1,]    1    0    0
## [2,]    0    0    1
## [3,]    0    1    0
```

## Complexité Théorique

La complexité de l'algorithme naïf est de  $O(n!)$ , où  $n$  est la taille de la matrice d'association. Cette complexité est due à la génération de toutes les permutations possibles des affectations et au calcul du coût pour chaque permutation.

### Génération de toutes les permutations :

- La fonction utilise la fonction `next_permutation`, qui génère toutes les permutations possibles des indices de 0 à  $n-1$ .
- Il y a  $n!$  permutations possibles dans le pire des cas, car il y a  $n$  éléments à permuter.
- La complexité de la génération de chaque permutation est  $O(n)$ , car elle nécessite un tri lexicographique des indices.

### Calcul du coût pour chaque permutation :

- Pour chaque permutation générée, le coût de cette permutation est calculé en utilisant la fonction `calculateCost`.
- La fonction `calculateCost` parcourt chaque ligne de la matrice d'association une fois, ce qui prend  $O(n)$ .
- Dans le pire des cas, la fonction `calculateCost` est appelée  $n!$  fois, une fois pour chaque permutation.

Ainsi, la complexité totale de l'algorithme naïf est de  $O(n!)$ , car la génération de toutes les permutations possibles et le calcul du coût pour chaque permutation contribuent tous deux à cette complexité. Cette approche est pratique pour de petites tailles de matrices, mais devient rapidement inefficace pour des tailles de problèmes plus importantes en raison de sa complexité.

## Complexité expérimentale (RCP)

```
temps_execution <- function(n_liste,p,algo) {
  #n_liste correspond aux dimensions des matrices de coûts générées,
  #p correspond au nombre de fois où l'on exécute l'algorithme pour ensuite moyenner le temps

  temps <- numeric(length(n_liste))

  for (j in seq_along(n_liste)) {
```

```

n <- n_liste[j]
cout <- matrice_couts(n)

start_time <- Sys.time()

for (i in 0:(p-1)) {
  algo(cout)
}

end_time <- Sys.time()
temps[j] <- (end_time - start_time)/p
}
return(temps)
}

```

Pour obtenir la complexité, on détermine la régression  $\log(T)=\log(n!)$ . Avoir un coefficient de régression de 1 signifierait que l'on a bien la même complexité que celle théorique.

```

liste_dimension <- seq(from = 6, to = 12, by = 1)

logn_factorielle <- log(factorial(liste_dimension))

temps_exe <- log(temps_execution(liste_dimension,5,NaiveAlgorithme))

reg <- lm(temps_exe~logn_factorielle)

summary(reg)

```

```

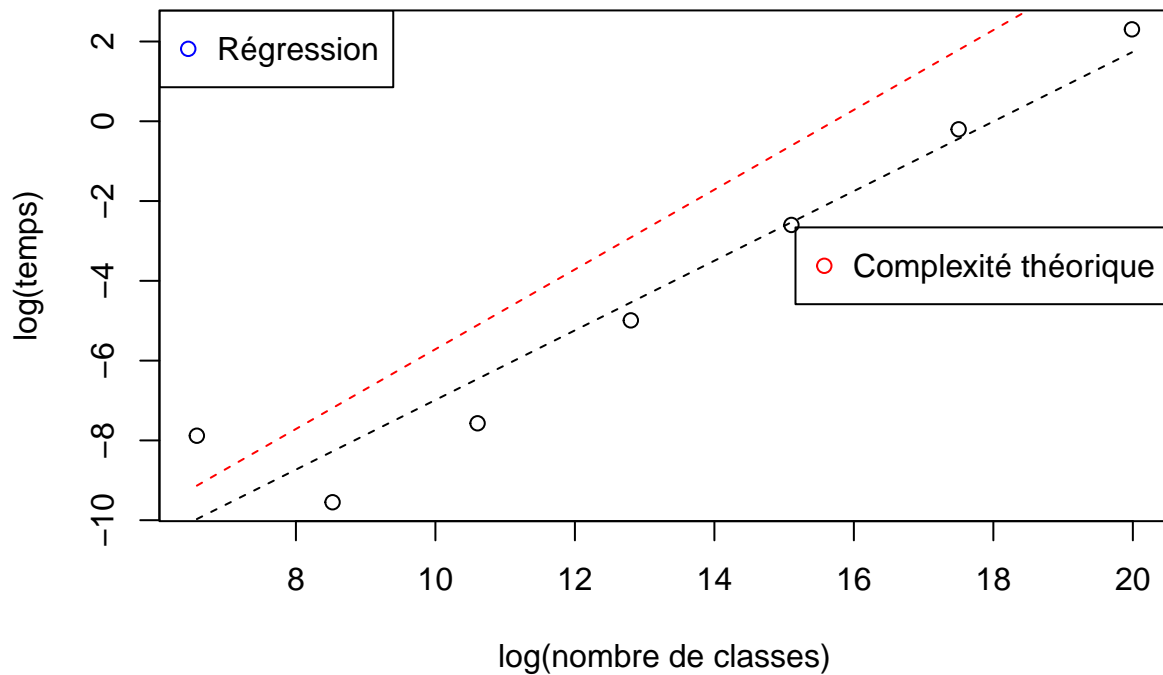
##
## Call:
## lm(formula = temps_exe ~ logn_factorielle)
##
## Residuals:
##      1      2      3      4      5      6      7
##  2.08952 -1.27603 -1.11164 -0.44930 -0.06638  0.23836  0.57547
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -15.7141     1.4524 -10.819  0.000117 ***
## logn_factorielle  0.8727     0.1055   8.271  0.000422 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.251 on 5 degrees of freedom
## Multiple R-squared:  0.9319, Adjusted R-squared:  0.9183
## F-statistic: 68.4 on 1 and 5 DF, p-value: 0.0004216

```

```

plot(logn_factorielle,temps_exe,xlab='log(nombre de classes)',ylab='log(temps)')
lines(logn_factorielle,coef(reg)[1]+logn_factorielle*coef(reg)[2],lty='dashed')
lines(logn_factorielle,coef(reg)[1]+logn_factorielle,col='red',lty='dashed')
legend("right", legend = "Complexité théorique", col = "red", pch = 1)
legend("topleft", legend = "Régression", col = "blue", pch = 1)

```



```
print(coef(reg)[2])
```

```
## logn_factorielle
##      0.8727075
```

La régression donne des valeurs très proches de 1 avec un  $R^2=0.99$ . Ainsi, on a bien une complexité en  $N!$ .

## Complexité expérimentale (R)

Vérifions maintenant que c'est aussi le cas pour l'algorithme en R:

```
temps_execution_R <- function(n_liste,p,algo) {
  #n_liste correspond aux dimensions des matrtrices de couts générées,
  #p correspond au nombre de fois où l'on exécute l'algorithme pour ensuite moyenner le temps

  temps <- numeric(length(n_liste))

  for (j in seq_along(n_liste)) {
    n <- n_liste[j]
    cout <- matrix(unlist(matrice_couts(n)), nrow = n, byrow = TRUE)

    start_time <- Sys.time()

    for (i in 0:(p-1)) {
```

```

    algo(cout)
  }

  end_time <- Sys.time()
  temps[j] <- (end_time - start_time)/p
}
return(temps)
}

```

```

liste_dimension <- seq(from = 2, to = 7, by = 1)

logn_factorielle <- log(factorial(liste_dimension))

temps_exe_R <- temps_execution_R(liste_dimension,5,NaiveAlgorithme_R)

log_temps_R <- log(temps_exe_R)

reg <- lm(log_temps_R~logn_factorielle)

summary(reg)

```

```

##
## Call:
## lm(formula = log_temps_R ~ logn_factorielle)
##
## Residuals:
##      1      2      3      4      5      6
## -0.02025  0.25743 -0.15701 -0.13895 -0.11094  0.16972
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -8.55151    0.14840  -57.62 5.43e-07 ***
## logn_factorielle  0.93002    0.02942   31.61 5.97e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1948 on 4 degrees of freedom
## Multiple R-squared:  0.996, Adjusted R-squared:  0.995
## F-statistic: 999.4 on 1 and 4 DF, p-value: 5.967e-06

```

```

print(coef(reg)[2])

```

```

## logn_factorielle
##      0.9300155

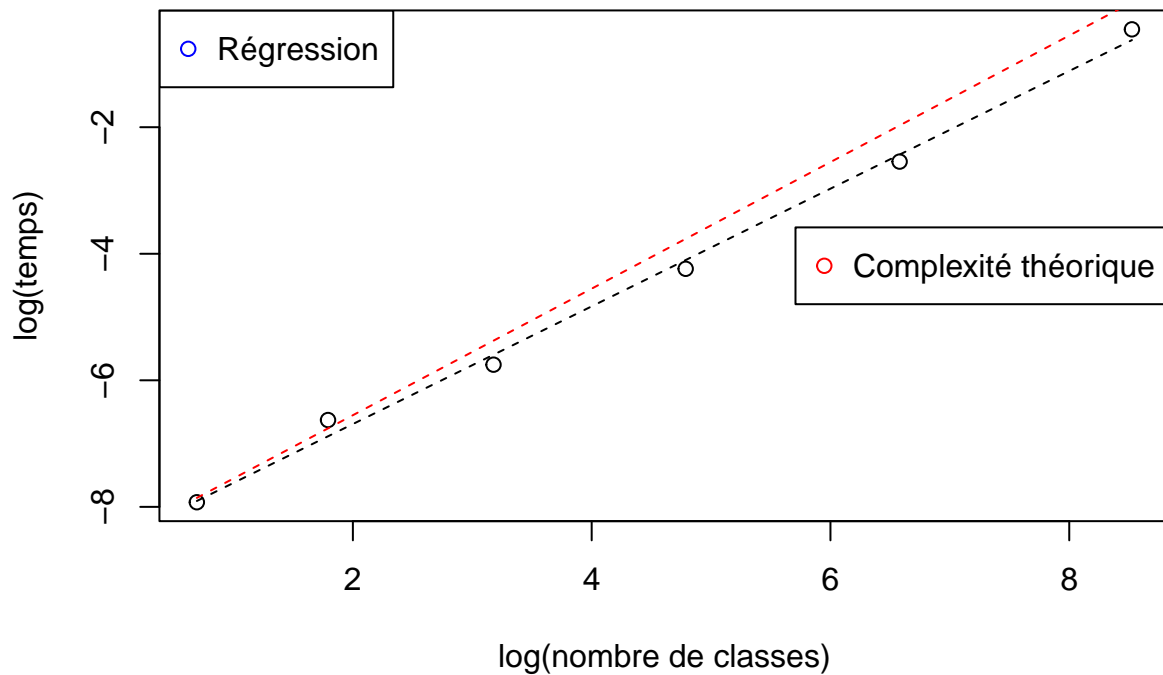
```

```

plot(logn_factorielle,log_temps_R,xlab='log(nombre de classes)',ylab='log(temps)')
lines(logn_factorielle,coef(reg)[1]+logn_factorielle*coef(reg)[2],lty='dashed')
lines(logn_factorielle,coef(reg)[1]+logn_factorielle,col='red',lty='dashed')
legend("right", legend = "Complexité théorique", col = "red", pch = 1)
legend("topleft", legend = "Régression", col = "blue", pch = 1)

```





A nouveau, le coefficient de régression est proche de 1 avec un  $R^2=0.98$ . La complexité expérimentale obtenu avec les deux algorithmes est donc similaire, et en accord avec la complexité théorique en  $N!$ .

## Vitesse R vs RCPP

```
# on a défini précédemment temps_exe_R <- temps_execution_R(liste_dimension,5,NaiveAlgorithme_R)
vitesse_cpp<-temps_execution(liste_dimension,5,NaiveAlgorithme)

print(vitesse_cpp)
```

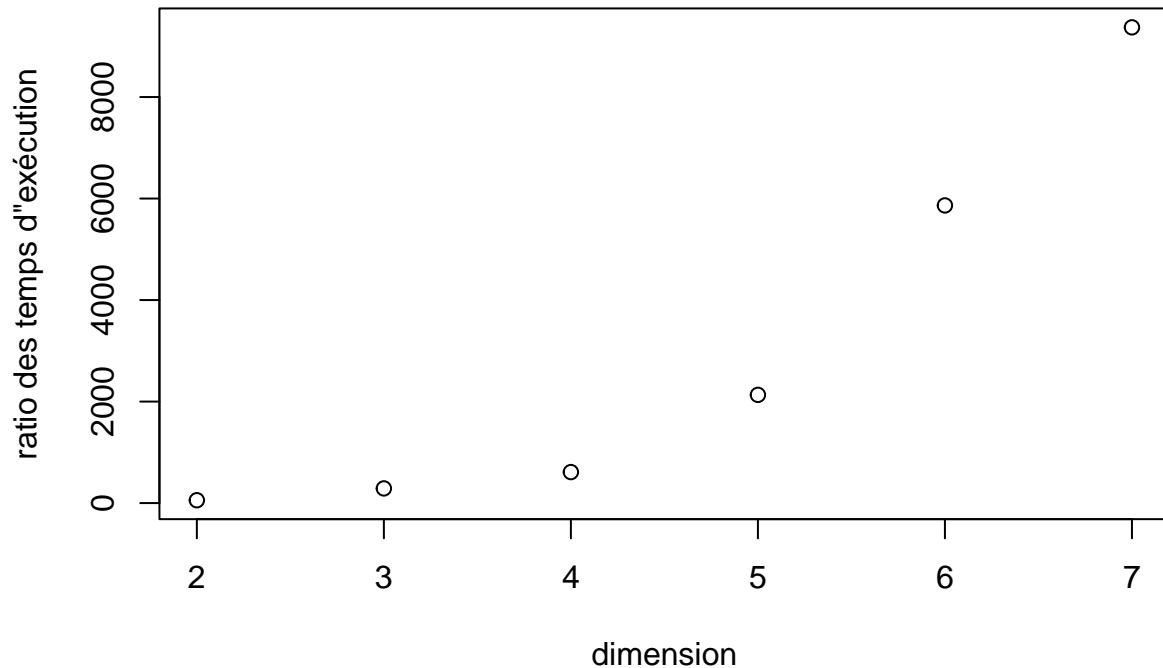
```
## [1] 6.437302e-06 4.577637e-06 5.197525e-06 6.771088e-06 1.339912e-05
## [6] 6.780624e-05
```

```
print(temps_exe_R)
```

```
## [1] 0.0003608227 0.0013231754 0.0031736374 0.0144361973 0.0785798073
## [6] 0.6355513573
```

```
plot(liste_dimension,temps_exe_R/vitesse_cpp,xlab='dimension',ylab='ratio des temps d"exécution',main='')
```

## temps d'exécution R/temps d'execution C++



On constate une différence de magnitude dans les temps d'exécution: comme attendu, le temps d'execution de c++ est beaucoup plus performant. A partir de 4 classes/clusters, le facteur entre les deux devient de l'ordre du millier, ce qui révèle l'utilité de C++ dans notre cas.

## Algorithme Hongrois

- Étape 1 : Soustraire le minimum de chaque ligne à tous les éléments de cette ligne, puis soustraire le minimum de chaque colonne à tous les éléments de cette colonne dans la matrice de coût.
- Étape 2 : Sélectionner le maximum de zéros indépendants, c'est-à-dire un seul zéro par ligne et par colonne, en parcourant tous les zéros non sélectionnés et en les sélectionnant s'ils ne partagent pas la même ligne ou colonne qu'un zéro déjà sélectionné.
- Étape 3-4-5 : Couvrir chaque colonne ayant un zéro sélectionné. Puis, pour chaque zéro non couvert, marquer la colonne de ce zéro et couvrir la ligne de ce zéro si aucune colonne n'est marquée pour ce zéro. Répéter jusqu'à ce qu'il n'y ait plus de zéros non couverts.
- Étape 6 : Trouver la valeur minimum des éléments non couverts dans la matrice. Ajouter cette valeur à toutes les lignes couvertes et la retirer à toutes les colonnes non couvertes. Recommencer à l'étape 1

L'itération de ces étapes permet de s'assurer que tout les chemins potentiellement optimaux sont explorés avant d'en choisir un. En effet, les étapes permettent de sélectionner un chemin et de s'assurer que celui-ci minimise la somme des coefficients sélectionnés par rapport au précédent puisque la matrice résultante à l'étape 6 est réutilisé pour le chemin suivant. L'algorithme sort de la boucle une fois qu'une solution optimale est trouvée, c'est-à-dire que réitérer les étapes n'a plus d'utilité.

```
sourceCpp("Hungarian.cpp")
source("Hungarian.R")
```

### Initialisation :

La fonction Hungarian est appelée, prenant en entrée une matrice de coûts matrix qui représente les coûts d'affectation entre deux ensembles d'éléments. Cette matrice peut être rectangulaire mais est transformée en une matrice carrée pour l'algorithme. Une copie de la matrice originale est faite pour sauvegarder les valeurs initiales.

La fonction adjust\_matrix est utilisée pour rendre la matrice d'entrée carrée en ajoutant des zéros aux rangées ou aux colonnes supplémentaires si nécessaire.

Des vecteurs temporaires sont initialisés pour suivre les affectations et les couvertures de lignes et de colonnes. La matrice M est initialisée pour stocker les affectations (zéros couverts et marqués).

### Fonctions auxiliaires :

Diverses fonctions auxiliaires sont utilisées pour effectuer des opérations spécifiques, telles que la recherche du plus petit élément dans une matrice, la recherche de zéros dans une matrice, la gestion des valeurs négatives, etc.

Exemple d'utilisation des packages:

- Package Rcpp:

```
cost_matrix <- list(
  c(9, 1, 5, 4),
  c(1, 4, 7, 3),
  c(8, 1, 6, 5),
  c(8, 1, 6, 5)
)
hungarian_adj_matrix <- Hungarian(cost_matrix, verbose = FALSE)

hungarian_adj_matrix
```

```
## [[1]]
## [1] 0 0 0 1
##
## [[2]]
## [1] 1 0 0 0
##
## [[3]]
## [1] 0 0 1 0
##
## [[4]]
## [1] 0 1 0 0
```

- Package R:

```

# Définir une matrice d'entrée
matrix_input <- matrix(c(9, 1, 8, 8, 1, 1, 4, 1, 5, 7, 6, 6, 4, 3, 5, 5), nrow = 4)

# Appliquer l'algorithme Hungarian
result <- Hungarian_R(matrix_input, verbose = FALSE)

# Afficher la solution trouvée
print(result)

```

```

##      [,1] [,2] [,3] [,4]
## [1,]    0    0    0    1
## [2,]    1    0    0    0
## [3,]    0    0    1    0
## [4,]    0    1    0    0

```

## Complexité Théorique

La complexité de l'algorithme hongrois est de  $O(n^3)$ , où  $n$  est la taille de la matrice carrée d'entrée. Cette complexité est due principalement à la recherche de chemins augmentants dans le graphe bipartite associé à la matrice d'entrée.

### Initialisation :

- Copier la matrice originale :  $O(n^2)$
- Ajuster la matrice en une matrice carrée :  $O(n^2)$
- Initialiser la matrice masquée  $M$ , les vecteurs de couverture des lignes et des colonnes, et d'autres variables temporaires :  $O(n^2)$

### Boucle principale :

- La boucle principale s'exécute jusqu'à ce que l'algorithme soit terminé, ce qui peut prendre plusieurs étapes.
- Chaque étape de la boucle principale implique généralement la recherche de chemins augmentants dans le graphe bipartite associé à la matrice d'entrée.
- La recherche de chemins augmentants utilise l'algorithme de l'arbre enraciné, qui a une complexité de  $O(n^3)$  dans le pire des cas.

### Calcul de la solution :

- Une fois que l'algorithme est terminé, la complexité de la dernière étape dépend de la manière dont la solution est extraite de la matrice d'affectation.
- Dans certains cas, cela peut impliquer une simple boucle sur la matrice d'affectation, ce qui ajoute une complexité supplémentaire de  $O(n^2)$ .

Ainsi, la complexité totale de l'algorithme hongrois est dominée par la recherche de chemins augmentants, ce qui donne une complexité globale de  $O(n^3)$ . Cette complexité fait de l'algorithme hongrois une méthode efficace pour résoudre le problème d'assignation pondérée dans de nombreuses applications.

On fonction de l'implementation de l'algorithme hongrois on peut avoir une autre complexité :

Les différentes étapes sont faites avec une complexité de 1 sur la totalité de la matrice de coût donc ont une complexité maximale de  $O(n^2)$ . De plus, Une itération de ces étapes garantit d'avoir au moins un zéro indépendant, ce qui indique que l'algorithme résout le problème au maximum en  $n^2$  itérations, correspondant au cas le plus extrême, à devoir itérer sur chacun des coefficients de la matrice d'assignation, qui a  $n^2$  coefficients. On en déduit que l'algorithme a une complexité théorique en  $O(n^4)$

## Complexité Expérimentale

```
liste_dimension <- seq(from = 10, to = 450, by = 20)

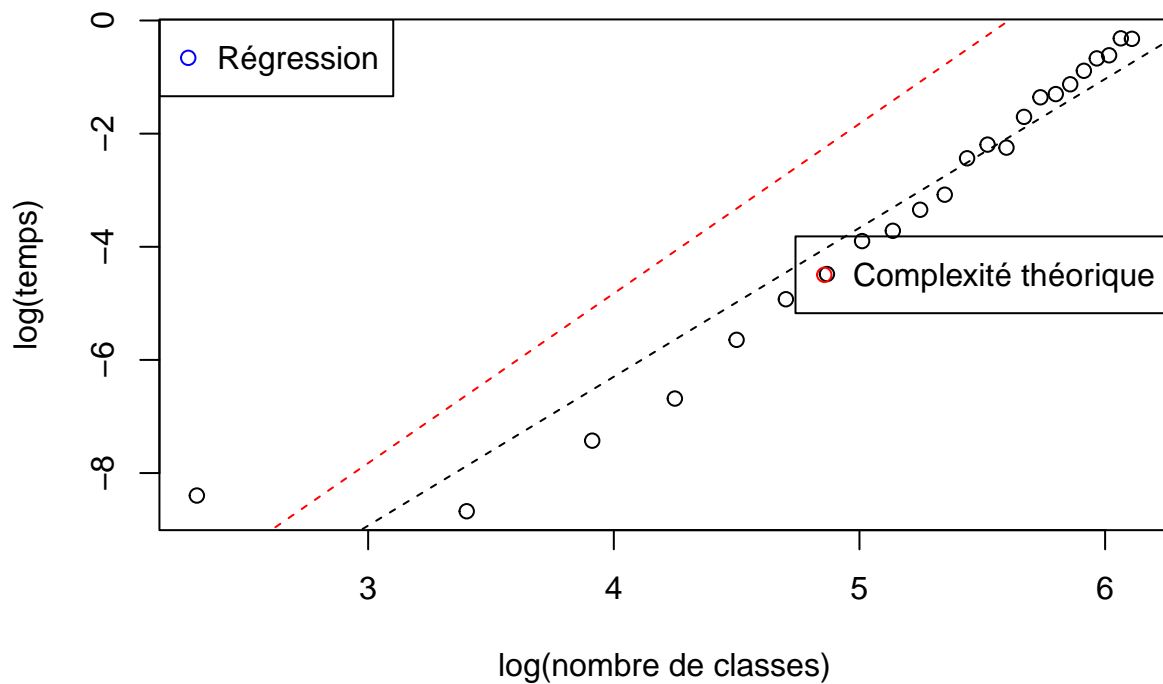
temps_exe_hong <- log(temps_execution(liste_dimension,5,Hungarian))

reg <- lm(temps_exe_hong~log(liste_dimension))

summary(reg)

##
## Call:
## lm(formula = temps_exe_hong ~ log(liste_dimension))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.03968 -0.43850  0.07692  0.37102  2.36689
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -16.8244     0.8387  -20.06 3.53e-15 ***
## log(liste_dimension)  2.6318     0.1604   16.41 1.88e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7171 on 21 degrees of freedom
## Multiple R-squared:  0.9277, Adjusted R-squared:  0.9242
## F-statistic: 269.4 on 1 and 21 DF,  p-value: 1.875e-13

plot(log(liste_dimension),temps_exe_hong,xlab='log(nombre de classes)',ylab='log(temps)')
lines(logn_factorielle,coef(reg)[1]+logn_factorielle*coef(reg)[2],lty='dashed')
lines(logn_factorielle,coef(reg)[1]+logn_factorielle*3,col='red',lty='dashed')
legend("right", legend = "Complexité théorique", col = "red", pch = 1)
legend("topleft", legend = "Régression", col = "blue", pch = 1)
```



```
print(coef(reg)[2])
```

```
## log(liste_dimension)
##          2.631779
```

La complexité expérimentale est là encore en accord avec la complexité théorique: nous avons bien un coefficient directeur proche de 3 avec un  $R^2$  solide, ce qui indique une complexité en  $N^3$ .

## Comparaison de vitesse R et C++

Nous générons une matrice de coût identique pour les deux algorithmes.

```
generateCostMatrix <- function(n, k) {
  matrix <- vector("list", n)
  for (i in 1:n) {
    matrix[[i]] <- vector("numeric", n)
    for (j in 1:n) {
      matrix[[i]][[j]] <- (i-1)*n + (j-1) + ((i-1) + (j-1)) %% k
    }
  }
  return(matrix)
}
```

```
ConvertMatrixListToMatrix <- function(matrixList){
  return(matrix(unlist(matrixList), nrow = length(matrixList), byrow = TRUE))
}
```

```
cost_matrix <- ConvertMatrixListToMatrix(generateCostMatrix(6, 4))
cat("Generated Matrix:\n")
```

```
## Generated Matrix:
```

```
print(cost_matrix)
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,]    0    2    4    6    4    6
## [2,]    7    9   11    9   11   13
## [3,]   14   16   14   16   18   20
## [4,]   21   19   21   23   25   23
## [5,]   24   26   28   30   28   30
## [6,]   31   33   35   33   35   37
```

Chaque algorithme est exécuté plusieurs fois afin d'obtenir une moyenne et un intervalle de confiance.

## Test de l'algorithme Rcpp

```
cost_matrix <- generateCostMatrix(100, 10)

# Nombre de fois à exécuter l'algorithme
num_executions <- 100

# Vecteur pour stocker les temps d'exécution
execution_times <- numeric(num_executions)

# Exécuter l'algorithme plusieurs fois et mesurer le temps d'exécution à chaque fois
for (i in 1:num_executions) {
  start_time <- Sys.time()
  # Appeler la fonction pour exécuter l'algorithme hongrois
  hungarian_adj_matrix <- Hungarian(cost_matrix)
  end_time <- Sys.time()
  execution_times[i] <- as.numeric(end_time - start_time, units = "secs") * 1000 # en millisecondes
}

# Calculer la moyenne des temps d'exécution
mean_execution_time <- mean(execution_times)

# Calculer l'écart type
standard_deviation <- sd(execution_times)

# Calculer l'intervalle de confiance à 95%
confidence_interval <- 1.96 * (standard_deviation / sqrt(num_executions))

# Afficher les résultats
cat("Moyenne du temps d'exécution :", mean_execution_time, "millisecondes\n")
```

```
## Moyenne du temps d'exécution : 0.2076674 millisecondes
```

```
cat("Intervalle de confiance (95%) : [", mean_execution_time - confidence_interval,
    ",", mean_execution_time + confidence_interval, "] millisecondes\n")
```

```
## Intervalle de confiance (95%) : [ 0.1787806 , 0.2365541 ] millisecondes
```

## Test de l'algorithme R

```
cost_matrix <- ConvertMatrixListToMatrix(generateCostMatrix(100, 10))

# Nombre de fois à exécuter l'algorithme
num_executions <- 100

# Vecteur pour stocker les temps d'exécution
execution_times <- numeric(num_executions)

# Exécuter l'algorithme plusieurs fois et mesurer le temps d'exécution à chaque fois
for (i in 1:num_executions) {
  start_time <- Sys.time()
  # Appeler la fonction pour exécuter l'algorithme hongrois
  hungarian_adj_matrix <- Hungarian_R(cost_matrix)
  end_time <- Sys.time()
  execution_times[i] <- as.numeric(end_time - start_time, units = "secs") * 1000 # en millisecondes
}

# Calculer la moyenne des temps d'exécution
mean_execution_time <- mean(execution_times)

# Calculer l'écart type
standard_deviation <- sd(execution_times)

# Calculer l'intervalle de confiance à 95%
confidence_interval <- 1.96 * (standard_deviation / sqrt(num_executions))

# Afficher les résultats
cat("Moyenne du temps d'exécution :", mean_execution_time, "millisecondes\n")
```

```
## Moyenne du temps d'exécution : 8.390167 millisecondes
```

```
cat("Intervalle de confiance (95%) :", mean_execution_time - confidence_interval,
    "à", mean_execution_time + confidence_interval, "millisecondes\n")
```

```
## Intervalle de confiance (95%) : 8.084446 à 8.695888 millisecondes
```

Pour l'algorithme hongrois en R, nous observons approximativement :

- Moyenne du temps d'exécution : 45.04352 millisecondes
- Intervalle de confiance (95%) : 43.7405 à 46.34653 millisecondes



Pour l'algorithme hongrois en Rcpp, nous observons approximativement :

- Moyenne du temps d'exécution : 0.2169228 millisecondes
- Intervalle de confiance (95%) : 0.1989332 à 0.2349123 millisecondes

Tandis que pour l'algorithme hongrois en C++, nous obtenons environ :

- Moyenne du temps d'exécution : 0.79411 millisecondes
- Intervalle de confiance (95%) : [0.343787, 1.24443] millisecondes

Ces résultats mettent en évidence l'impact significatif de l'implémentation sur les performances des algorithmes, montrant clairement l'avantage de l'implémentation en C++ en termes de vitesse d'exécution.

## Performances sur MNIST

On utilise les clusters associés aux classes pour prédire le test de mnist. Pour ce faire, nous utilisons la matrice d'assignation classe/cluster obtenue pour assigner les classes aux clusters prédits par l'algorithme des k-means. Nous comparons ensuite avec la vraie classe des données test mnist avec celle prédite par cette association.

```
# matrice bijection classes/clusters trouvés par l'algorithme hongrois
assignation <- Hungarian(cost_matrix_mnist)

# version liste de la bijection
association <- numeric(k)

for (i in seq_along(assignation)) {

  association[i] <- which(assignation[[i]] == 1)
}

# assignement des données à des classes prédites par la bijection
assignement_classe <- association[assignement_cluster]

précision <- mean(assignement_classe != y_train) * 100

print(précision)
```

```
## [1] 91.21
```

Nous obtenons une efficacité supérieure à 95%. Nous pouvons donc supposer au-delà du fait que la structure des classes est bien captée par l'algorithme kmeans que l'algorithme hongrois permet de trouver une association cohérente entre les clusters et les classes. On peut donc potentiellement se reposer sur celui-ci pour labeliser de nouvelles données après clusterisations.

## Conclusion

Nous avons pu appliquer l'algorithme hongrois à l'association de classes de données et de clusters prédits par kmean. pour ceci, nous avons défini une matrice de coût basée sur la distance euclidienne entre centroïdes

des classes et des clusters puis nous avons appliqué un algorithme naïf de résolution de l'association qui s'est avéré très demandeur en ressources du fait de sa complexité factorielle. Nous avons donc appliqué l'algorithme hongrois qui réduit la complexité à une puissance polynomiale de la dimension et avons pu remarquer qu'utiliser RCPP était largement avantageux en terme de rapidité. Enfin, une comparaison entre les associations obtenues et les classes réelles ont montrées que l'algorithme hongrois permettait effectivement de trouver une solution optimale.