

# Segunda entrega de projeto

Karen Kaori Yonea - 10349471

## 1 Projeto4

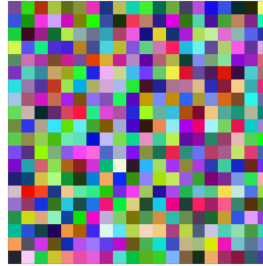
### 1.1 Implementação

```
1 SOM <- function(train, n) {
2   train <- train
3   n <- n
4   n_train <- dim(train)[1]
5   grid <- meshgrid(1:n, 1:n)
6   c <- matrix(c(as.vector(t(grid$X)), as.vector(t(
7     grid$Y))), ncol = 2)
8
9   #Definir a saída
10  res <- list()
11
12  #Configuração inicial: pesos aleatórios
13  nodes <- rand(n*n, dim(train)[2])
14  res$init = nodes
15
16  #Pesos atualizados
17  res$trained <- training(c)
18
19  #Retorna a matriz de pesos inicial e final
20  return(res)
21
22  training <- function(c) {
23    for (i in 1:dim(train)[1]) {
24      #Encontrar o neurônio "vencedor"
25      match = matching(train[i,])
26      #Distâncias
27      d <- rowNorms(c - c[match], method='
28        euclidean')
29      L = learning_rate(i)
30      S = learning_radius(i)
31      aux = train[i,] - nodes
32      #Atualizar pesos
33      nodes = nodes + L*(S * aux)
34    }
35    return(nodes)
36  }
37
38  matching <- function(teachers) {
39    norms = rowNorms(nodes - teachers, method='
40      euclidean')
41    match = which.min(norms)
42    return(match)
43  }
44
45  neighbourhood <- function(t) {
46    hfl = n_train/4
47    ini = n/2
48    nei = ini*exp(-t/hfl)
49    return(nei)
50  }
```

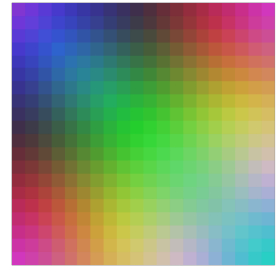
```
50 learning_rate <- function(t) {
51   hfl = n_train/4
52   ini = 0.1
53   rat = ini*exp(-t/hfl)
54   return(rat)
55 }
56
57 learning_radius <- function(t) {
58   s = neighbourhood(t)
59   rad = exp(-d^2/(2*s^2))
60   return(rad)
61 }
```

### 1.2 Paleta de cores

Para testar o algoritmo foi utilizado dados gerados de forma aleatório em escala de cor RGB:



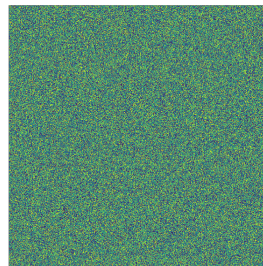
(a) Antes de aplicar o SOM



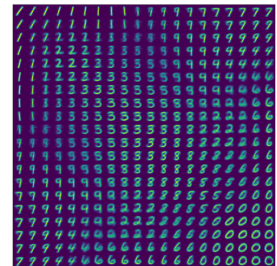
(b) Depois de aplicar o SOM

### 1.3 Base MNIST(caracteres manuscritos)

Utilizando a base de dados de imagens de teste do MNIST com 10 categorias obteve-se:



(a) Antes de aplicar o SOM



(b) Depois de aplicar o SOM

## 2 Projeto 5

### 2.1 Implementação

#### 2.1.1 KNN

```
1 #Análise da base de treino
2 guess.knn <- function(x, train, trainlabels, k) {
3   #Distância euclidiana
4   xmatrix <- matrix(as.numeric(x), nrow=nrow(
5     train), ncol=length(x), byrow=T)
6   xmatrix <- (abs(as.matrix(train)-xmatrix))^2
7   diffs <- (rowSums(xmatrix))^(1/2)
8   #Ordena as distâncias
9   diffs <- data.frame(dist=diffs,label=
10     trainlabels)
11   diffs <- (diffs[order(diffs$dist),])
12   diffs <- diffs[1:k,]
13
14   guess <- names(sort(-table(diffs$label)))[1]
15
16   return(guess)
17 }
18
19 knn <- function(train, test, trainlabels, k) {
20   #Formatando a base de treino
21   subsample <- sample(1:nrow(train), nrow(train),
22     replace=F)
23   train <- train[subsample,]
24   trainlabels <- trainlabels[subsample]
25
26   #Atribuir os valores das classes da base teste
27   kn <- apply(test, 1, function(x) guess.knn(x,
28     train, trainlabels, k))
29
30   return(kn)
31 }
```

#### 2.1.2 K-means

```
1 k_means <- function(dataset, k, threshold=1e-3) {
2   #Definindo os centroides iniciais
3   ids = sample(1:nrow(dataset), size=k)
4   centroid = as.data.frame(dataset[ids,])
5   div = 2 * threshold
6
7   #Divergente > que o limiar
8   while (div > threshold) {
9     # Para cada centroide
10     dists = NULL
11     for (i in 1:k) {
12       # calcular a distancia euclidiana de cada
13       # ponto do
14       # espaço em relação a cada centroide
15       euclidean = apply(dataset, 1, function(row) {
16         {sqrt(sum((row - centroid[i,])^2))})
17       }
18       dists = cbind(dists, euclidean)
19     }
20
21     # Descobrir qual centroide cada ponto está
22     # associado
23     clusters = apply(dists, 1, function(row) {
24       which.min(row) })
25
26     # Atualizar a posição dos centroides
27     div = 0
28     for (i in 1:k) {
29       ids = which(clusters == i)
30       position = colMeans(as.data.frame(dataset[
31         ids,]))
32       div = div + sqrt(sum((centroid[i,] -
33         position)^2))
34       centroid[i,] = position
35     }
36
37     div = div / k
38   }
39
40   km = list()
41   km$k = k
42   km$centroid = centroid
43   km$cluster = clusters
44
45   return(km)
46 }
```

### 2.2 Sinal 1D - Autômatos probabilísticos

Para essa parte foi gerado 500 padrões de tamanho  $M = 500$  para cada autômato, tendo um total de 1500 dados.

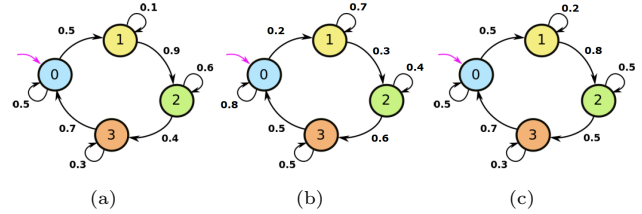


Figura 3: Autômatos utilizados

A escolha de dois autômatos muito similares para testar os métodos foi proposital.

Atributos escolhidos:

- Média do tamanho dos bursts
- Entropia do tamanho dos bursts
- Média das distâncias intersímbolos
- Entropia das distâncias intersímbolos

#### 2.2.1 KNN

Foram selecionados aleatoriamente 250 índices para treino e outros 250 para teste.

Matriz de Confusão

|   | a  | b  | c  |
|---|----|----|----|
| a | 82 | 0  | 13 |
| b | 0  | 82 | 0  |
| c | 12 | 0  | 61 |

#### 2.2.2 K-means

Utilizando 250 índices foram obtidos resultados pouco assertivos, então foi tomada a escolha de analisar 500 valores.

Matriz de Confusão

|   | a   | b   | c   |
|---|-----|-----|-----|
| a | 134 | 0   | 26  |
| b | 1   | 164 | 0   |
| c | 11  | 0   | 164 |

Os melhores resultados foram obtidos utilizando entre (500, 750) índices.

### 2.3 Base MNIST(caracteres manuscritos)

#### 2.3.1 KNN

Para essa parte foram selecionados aleatoriamente 250 índices da base de treino e 250 da base de teste.

Matriz de Confusão

|   | 0  | 1  | 2  | 3  | 4  | 5 | 6  | 7  | 8 | 9  |
|---|----|----|----|----|----|---|----|----|---|----|
| 0 | 20 | 0  | 0  | 0  | 0  | 0 | 0  | 0  | 0 | 0  |
| 1 | 0  | 32 | 0  | 0  | 0  | 0 | 0  | 0  | 2 | 0  |
| 2 | 2  | 8  | 12 | 0  | 2  | 0 | 0  | 1  | 1 | 0  |
| 3 | 0  | 3  | 0  | 15 | 0  | 1 | 0  | 0  | 0 | 3  |
| 4 | 0  | 4  | 0  | 0  | 21 | 0 | 1  | 0  | 0 | 0  |
| 5 | 3  | 1  | 0  | 5  | 0  | 9 | 0  | 1  | 1 | 1  |
| 6 | 0  | 0  | 0  | 0  | 2  | 0 | 20 | 0  | 0 | 0  |
| 7 | 0  | 3  | 0  | 0  | 2  | 0 | 0  | 19 | 0 | 7  |
| 8 | 1  | 3  | 0  | 7  | 1  | 0 | 0  | 1  | 5 | 3  |
| 9 | 0  | 2  | 0  | 0  | 6  | 0 | 0  | 0  | 1 | 21 |

### 2.3.2 K-means

Para essa parte foram selecionados aleatoriamente 250 índices da base de teste e os atributos foram pré-processados com a função `nearZeroVar` que diagnostica preditores que possuem um valor único (ou seja, preditores de variação zero) ou preditores que possuem as duas características a seguir: eles têm muito poucos valores exclusivos em relação ao número de amostras e a razão da frequência do valor mais comum para a frequência do segundo valor mais comum é grande.

Matriz de Confusão

|   | 0  | 1  | 2  | 3 | 4  | 5 | 6  | 7  | 8 | 9  |
|---|----|----|----|---|----|---|----|----|---|----|
| 0 | 26 | 0  | 0  | 0 | 0  | 1 | 0  | 1  | 0 | 0  |
| 1 | 0  | 27 | 0  | 0 | 0  | 0 | 0  | 0  | 2 | 0  |
| 2 | 0  | 7  | 12 | 0 | 6  | 0 | 0  | 0  | 1 | 0  |
| 3 | 0  | 5  | 0  | 9 | 1  | 2 | 0  | 1  | 2 | 1  |
| 4 | 0  | 0  | 0  | 0 | 12 | 1 | 0  | 1  | 5 | 11 |
| 5 | 0  | 1  | 0  | 8 | 1  | 4 | 0  | 0  | 5 | 2  |
| 6 | 0  | 0  | 0  | 0 | 8  | 1 | 12 | 0  | 1 | 0  |
| 7 | 1  | 1  | 0  | 0 | 0  | 0 | 0  | 16 | 0 | 5  |
| 8 | 0  | 5  | 0  | 2 | 0  | 6 | 0  | 1  | 9 | 1  |
| 9 | 0  | 0  | 0  | 0 | 4  | 2 | 0  | 2  | 0 | 18 |

Alguns grupos ficam bem definidos como o zero enquanto outros ficam indistinguíveis

## 2.4 Base IRIS

### 2.4.1 KNN

Foram selecionados aleatoriamente 75 índices da base de dados como treino e os outros 75 índices como teste.

Matriz de Confusão

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 25     | 0          | 0         |
| versicolor | 0      | 25         | 0         |
| virginica  | 0      | 5          | 20        |

### 2.4.2 K-means

Utilizou-se todas as entradas da base de dados e não foi feito nenhum tipo de pré-processamento.

Matriz de Confusão

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 50     | 0          | 0         |
| versicolor | 0      | 48         | 2         |
| virginica  | 0      | 14         | 36        |

## 3 Projeto 6

### 3.1 Implementação

PCA:

```
1 pca <- function(p) {  
2   k = cov(p)  
3   Q = eigen(k)  
4   p = scale(p)  
5   for (i in 1:dim(p)[1]) {  
6     x = t(p[i,]) %*% Q$eigenvectors  
7     p[i,] = x  
8   }  
9   return(p)  
10 }
```

Método descrito por Newman:

```
1 spectral_clustering <- function(dt, k, sig) {  
2   #Definindo variáveis  
3   S <- dt  
4   n <- nrow(dt)  
5   W <- matrix(rep(0,n^2), nrow = n, ncol=n)  
6   D <- diag(n)  
7  
8   #Calculando matriz W  
9   for (i in 1:n){  
10    for (j in 1:n){  
11      if (i != j){  
12        W[i,j] <- exp( - sqrt(sum((S[i,]-S[j,])  
13          ^2)) / 2*sig)  
14      }  
15    }  
16  }  
17  
18  #Calculando matriz diagonal  
19  for (i in 1:n){  
20    D[i,i] <- sum (W[i,])  
21  }  
22  
23  #Definindo L  
24  aux = solve(sqrt(D))  
25  L <- aux %*% W %*% aux  
26  
27  #Autovetores de L (funcao eigen ja ordena)  
28  eig <- eigen(L)$vectors  
29  X <- eig[, (1 : k)]  
30  
31  #Definindo matriz de comunidade  
32  Y <- matrix(0,nrow=n,ncol=k)  
33  com <- matrix(0,nrow=n,ncol=1)  
34  
35  #Encontra as comunidades  
36  for (i in 1:n){  
37    for (j in 1:k){  
38      Y[i,j] <- X[i,j] / (sqrt (sum(X[i,])^2))  
39    }  
40  }  
41  
42  return(Y)
```

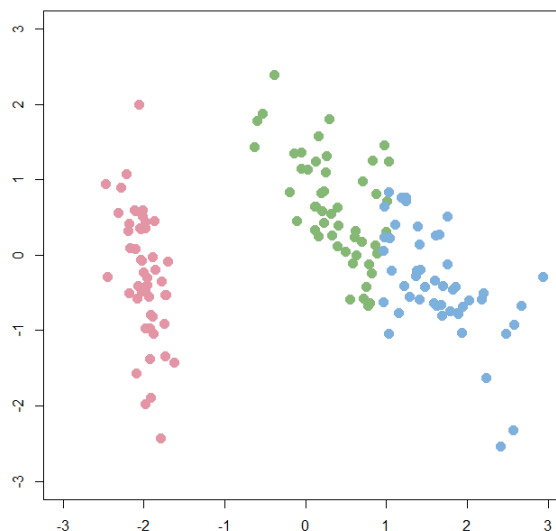


Figura 4: PCA com a separação das comunidades

### 3.2 Resultados

Matriz de Confusão Método de Newman

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 50     | 0          | 0         |
| versicolor | 0      | 44         | 6         |
| virginica  | 0      | 6          | 44        |

Matriz de Confusão K-means

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 50     | 0          | 0         |
| versicolor | 0      | 48         | 2         |
| virginica  | 0      | 14         | 36        |

No geral ambos os métodos diferenciaram os grupos de forma satisfatória, porém, o custo computacional da implementação do k-means para essa precisão é muito maior, sendo assim, compensa nesse caso utilizar o método descrito por Newman.

## 4 Projeto 7

No projeto foram utilizadas as seguintes bibliotecas disponibilizadas do repositório CRAN:

- cluster: algoritmos de clusterização
- dendextend: comparar dois dendrogramas

### 4.1 Dendrogramas

Todos os dendrogramas foram gerados à partir da base de dados IRIS e utilizando a distância euclidiana entre os pontos.

#### 4.1.1 Single-linkage

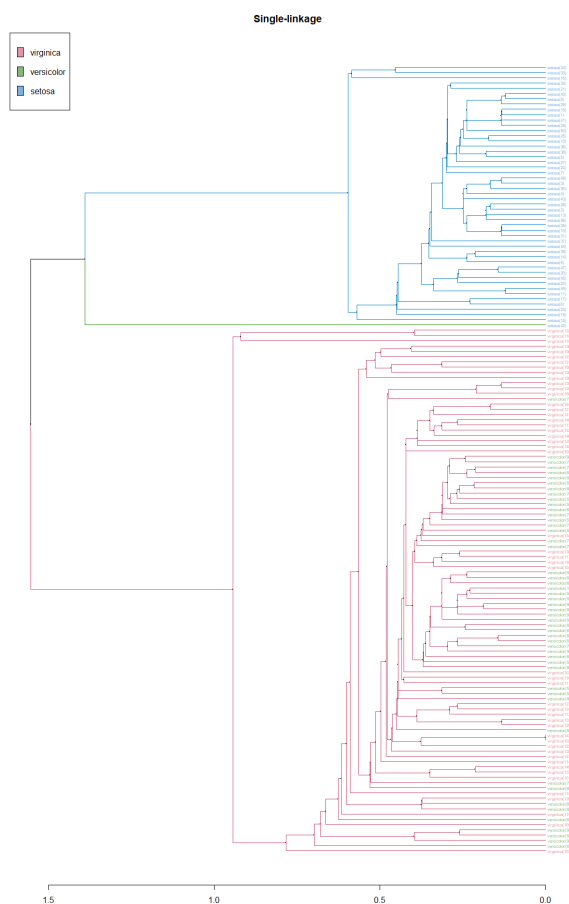


Figura 5: Dendrograma single-linkage

Matriz de Confusão

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 49     | 1          | 0         |
| versicolor | 0      | 0          | 50        |
| virginica  | 0      | 0          | 50        |

#### 4.1.2 Complete-linkage

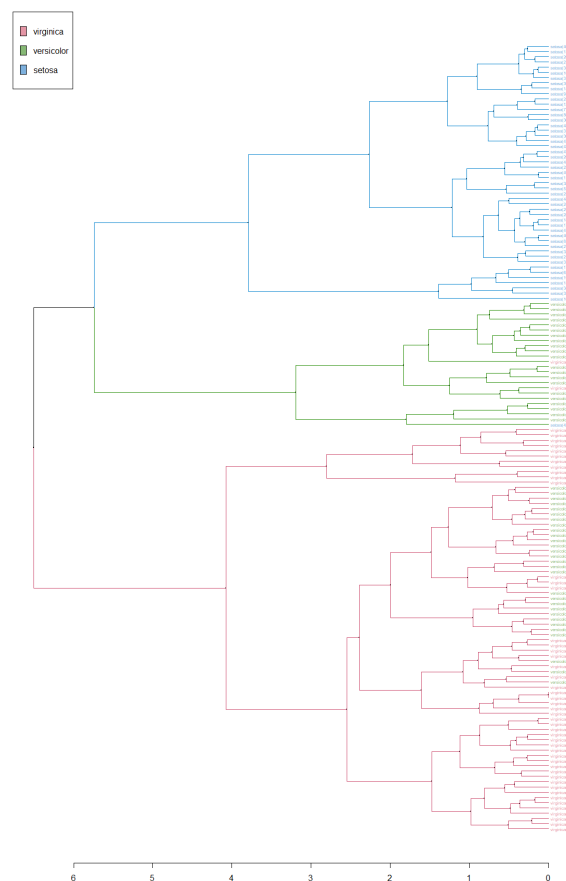


Figura 6: Dendrograma complete-linkage

Matrix de Confusão

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 49     | 1          | 0         |
| versicolor | 0      | 21         | 29        |
| virginica  | 0      | 2          | 48        |

### 4.1.3 Median

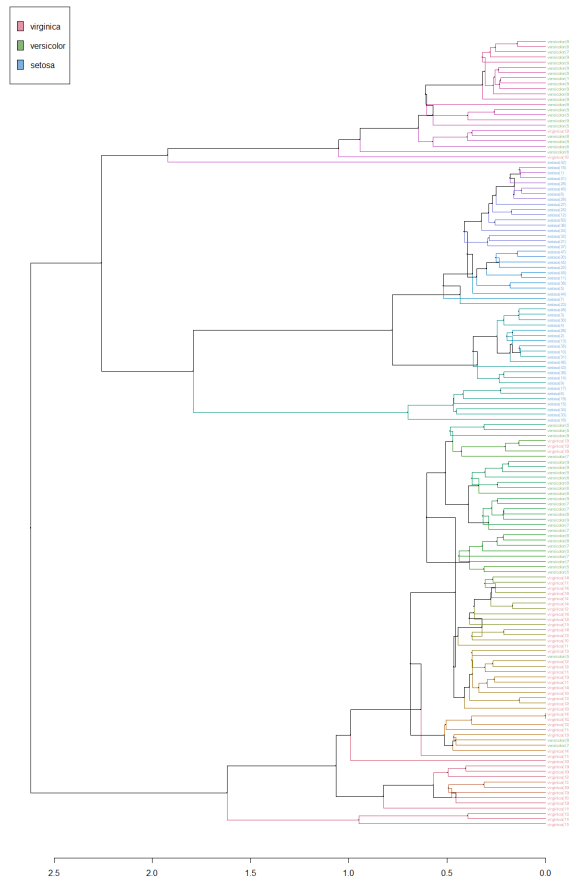


Figura 7: Dendrograma median

Matriz de Confusão

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 49     | 1          | 0         |
| versicolor | 0      | 21         | 29        |
| virginica  | 0      | 2          | 48        |

### 4.1.4 Ward

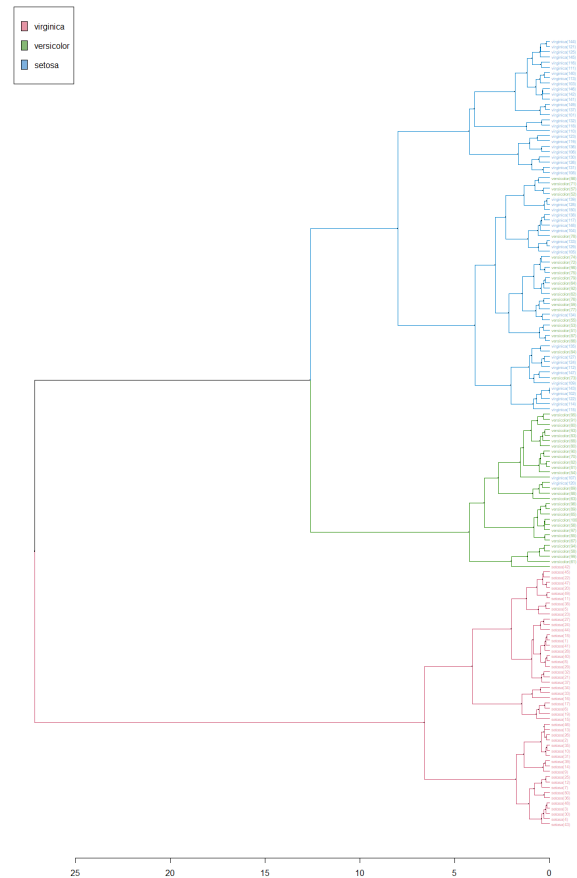


Figura 8: Dendrograma ward.D2

Matriz de Confusão

|            | setosa | versicolor | virginica |
|------------|--------|------------|-----------|
| setosa     | 49     | 1          | 0         |
| versicolor | 0      | 27         | 23        |
| virginica  | 0      | 2          | 48        |

## 4.2 Relação entre os dendrogramas

Para mensurarmos a relação entre os dendrogramas gerados pelos diferentes métodos, calcularemos o emaranhamento(entanglement) entre eles, essa medida assume valores entre 0 e 1, quanto menor, mais alinhados são os dendrogramas.

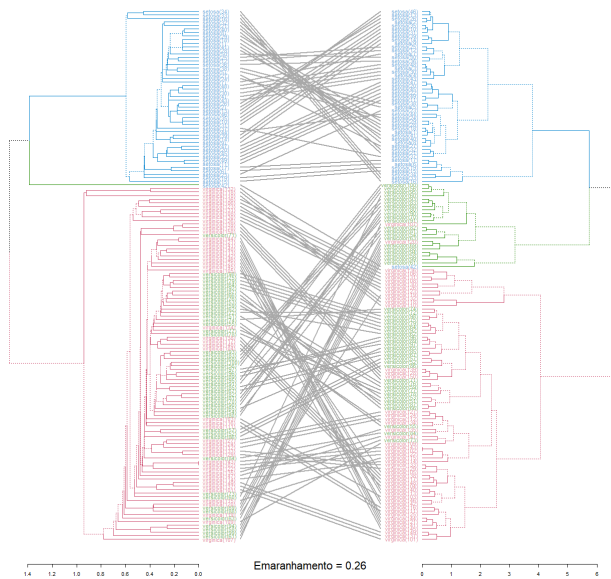


Figura 9: single e complete

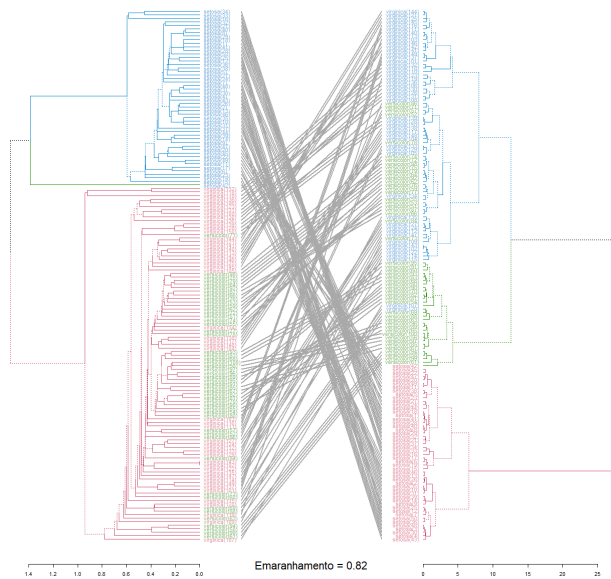


Figura 11: single e ward

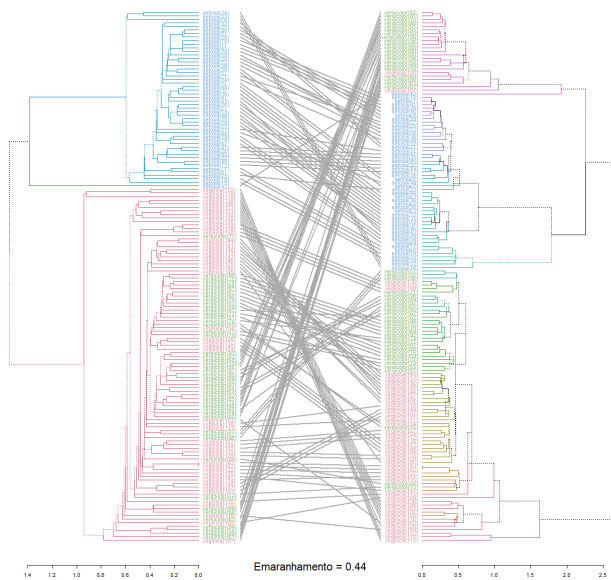


Figura 10: single e median

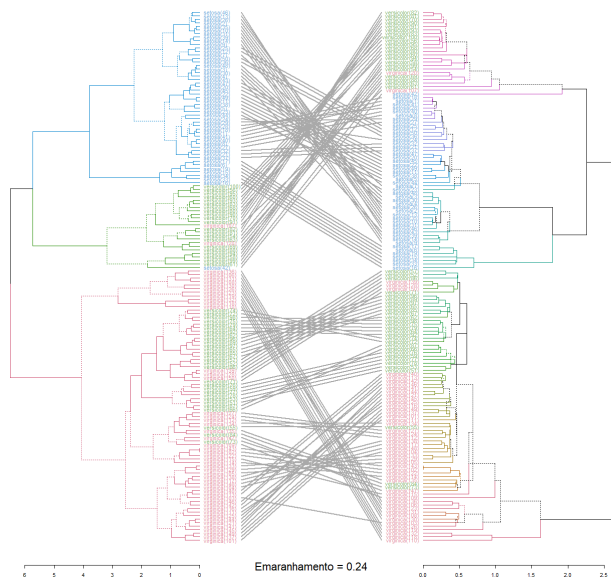


Figura 12: Complete e median

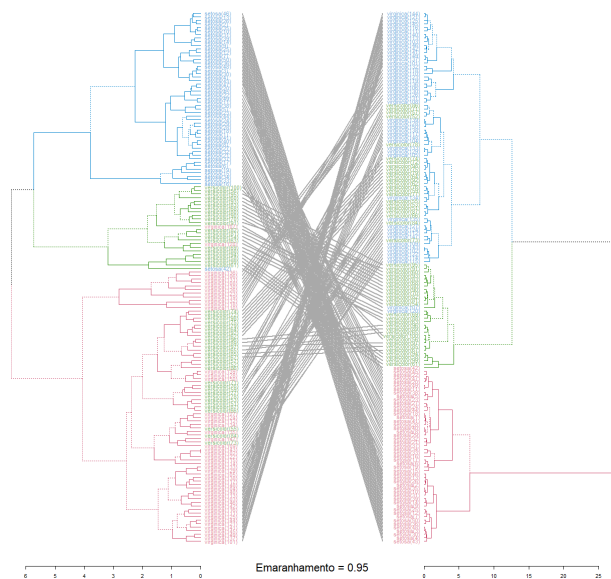


Figura 13: Complete e ward

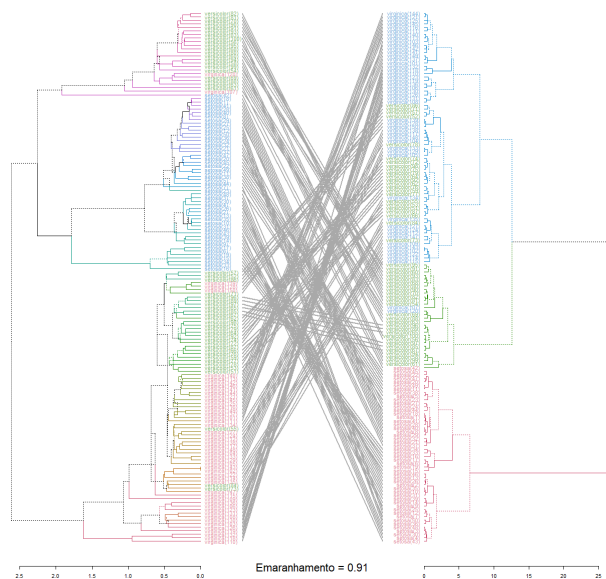


Figura 14: Median e ward

Ao obter todos os valores de emaranhamento par a par dos 4 métodos é tirada a conclusão de que os métodos: single-linkage, complete-linkage e median; agrupam os itens de forma parecida e completamente diferente da forma que o método ward.