

Evaluation des algorithmes de ranking pour classer des produits sur un site de e-commerce

Flavie Bertrand et Marion Tremblay

2025-12-08

1. Description du problème et objectif

1.1. Problème général de ranking

On considère :

- $i \in \{1, 2, \dots, n\}$: l'élément i
- un score associé à chaque élément : $s_i \in \mathbb{R}$
- une variable binaire $x_{i,p} \in \{0, 1\}$ qui vaut 1 si l'élément i est placé à la position p (0 sinon)
- un poids w_p associé à la position p

L'objectif est de déterminer un classement des k premiers éléments maximisant le score total.

$$\max_x \sum_{i=1}^n \sum_{p=1}^k w_p s_i x_{i,p}$$

Sous contraintes :

- Un élément ne peut occuper au plus qu'une seule position :

$$\sum_{p=1}^k x_{i,p} \leq 1, \quad \forall i$$

- Chaque position doit être occupée par exactement un élément :

$$\sum_{i=1}^n x_{i,p} = 1, \quad \forall p$$

- Contraintes de groupes :

On définit des groupes :

$$G_g \subseteq \{1, \dots, n\}, \quad g = 1, \dots, G$$

On impose que chaque groupe G_g apparaisse au plus α_g fois dans les k premiers :

$$\sum_{i \in G_g} \sum_{p=1}^k x_{i,p} \leq \alpha_g, \quad \forall g$$

1.2. Exemple concret

1.2.1. Contexte

On applique ce problème pour sélectionner les 5 produits les plus pertinents pour les afficher sur la page d'accueil d'un site de e-commerce. Chaque produit a un score s_i correspondant à la moyenne des notes données par les utilisateurs. On souhaite maximiser ce score tout en respectant certaines contraintes :

- Diversité des catégories : pas plus de 2 produits de la même catégorie
- Contrainte marketing : au moins un produit sponsorisé doit apparaître dans le top 5

Pour répondre à ce problème on utilise la base de données des retours produits du site Amazon datant de 2018. Nous avons sélectionné un échantillon de 250 produits de cette base de données. Nous avons généré une colonne en plus pour réaliser la contrainte marketing, en associant 15% des produits de la base de données à des produits sponsorisés.

Asin	Score	Category	Sponsored
B00004ZCJI	4,4	Electronics	False
B00004ZCJJ	4,4	Electronics	False
B00079ULA8	4,1	Sports	False
B000CBSNRY	4,9	Toys_Games	False
B0001YR54E	4,5	Clothing	True

Id	Score	Groupe
1	4,4	Electronics, not_sponsored
2	4,4	Electronics, not_sponsored
3	4,1	Sports, not_sponsored
4	4,9	Toys_Games, not_sponsored
5	4,5	Clothing, sponsored

1.2.2. Formulation mathématique

La fonction objectif est :

$$\max_x \sum_{i=1}^n \sum_{p=1}^5 w_p s_i x_{i,p}$$

où :

- $x_{i,p}$ est une variable binaire qui vaut 1 si l'élément i est placé à la position p (0 sinon)
- w_p est le poids associé à la position p : il vaut 5 à la position puis est décrémenter de 1 à chaque position suivante.

Sous contraintes

- Chaque position contient exactement un produit :

$$\sum_{i=1}^n x_{i,p} = 1 \quad \forall p$$

- Un produit ne peut être affiché qu'une seule fois :

$$\sum_{p=1}^5 x_{i,p} \leq 1 \quad \forall i$$

- Maximum 2 produits par catégorie :

$$\sum_{i \in \text{cat}} \sum_{p=1}^5 x_{i,p} \leq 2$$

- Au moins 1 produit sponsorisé dans le top 5 :

$$\sum_{i \in \text{sponsored}} \sum_{p=1}^5 x_{i,p} \geq 1$$

Cette dernière contrainte est équivalente à avoir maximum 4 produits non sponsorisé dans le top 5, soit :

$$\sum_{i \in \text{not sponsored}} \sum_{p=1}^5 x_{i,p} \leq 4$$

2. Approche heuristique

L'approche **heuristique** constitue une première stratégie simple pour résoudre le problème de ranking pondéré avec contraintes de groupes.

Principe algorithmique

1. Initialiser les compteurs de chaque groupe à 0 et liste des éléments sélectionnés vide.
2. Pour chaque position ($p = 1, \dots, k$) :
 - Considérer les éléments non encore sélectionnés respectant les contraintes de groupes.
 - Choisir l'élément avec le gain immédiat maximal $w_p \cdot s_i$.
 - Mettre à jour les compteurs de groupes et le score total.
3. Arrêt si aucune sélection possible pour une position ou si toutes les positions sont remplies.

Cette méthode est très rapide et facile à implémenter, mais elle ne garantit pas d'optimalité globale, car elle ne considère pas les effets futurs de chaque choix.

Exemple illustratif

On applique l'algorithme sur le jeu de données décrit en introduction :

```
## [1] 72.9
```

##	id	score	groupes	position	poids_position	score_pondere
## 1	40	5.0	Beauty,Not_sponsored	1	5	25.0
## 2	4	4.8	Electronics,Not_sponsored	2	4	19.2
## 3	8	4.8	Toys_Games,Not_sponsored	3	3	14.4
## 4	42	4.8	Automotive,Not_sponsored	4	2	9.6
## 5	5	4.7	Pet_Supplies,Sponsored	5	1	4.7

```
## [1] 72.9
```

- L'algorithme privilégie les éléments les mieux notés pour les premières positions.
- Les contraintes de groupe peuvent limiter le choix, ce qui peut conduire à des choix sous-optimaux pour les positions suivantes.
- Cette approche est rapide et simple.

3. Solution 1 : programmation dynamique

L'algorithme

Pour résoudre le problème de ranking sous contraintes de groupe et de poids positionnels, on utilise une programmation dynamique :

Chaque état est défini par :

- le nombre de positions déjà remplies,
- le nombre d'éléments choisis par groupe (compteurs de groupes).
- À chaque étape (position p), on essaie tous les éléments possibles qui respectent les contraintes de groupes et qui n'ont pas déjà été choisis.

À la fin, on récupère l'état final avec le score total maximal et la sélection correspondante.

La relation de récurrence est :

$$DP[s][c][S] = \max\{DP[p-1][c'][S'] + w[p]s[i]\}$$

avec

- p : nombre d'éléments
- $c = (c[1], c[2], \dots, c[G])$: vecteur des compteurs par groupe
- S : ensemble des éléments déjà utilisés
- $DP[s][c][S]$: meilleur score pour avoir sélectionné p éléments avec configuration c et ensemble S .

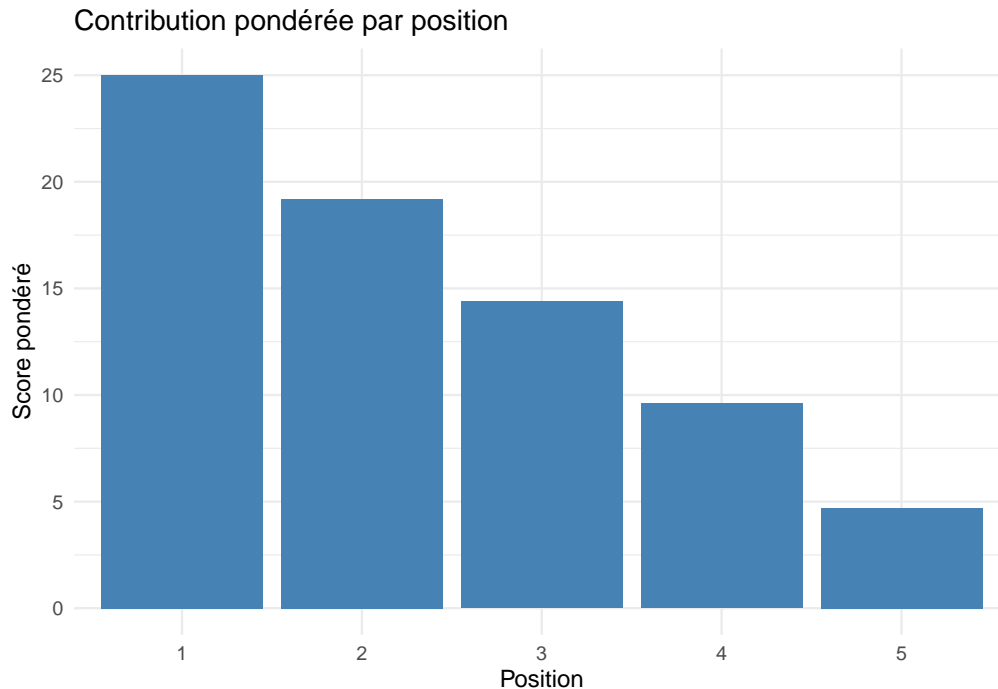
Exemple concret

On applique l'algorithme une nouvelle fois sur le jeu de données décrit en introduction :

Table 3: Tableau des éléments sélectionnés par l'algorithme ranking max

	Position	ID	Score	Groupe	Poids position	Score pondéré
40	1	40	5.0	Beauty , Not_sponsored	5	25.0
4	2	4	4.8	Electronics , Not_sponsored	4	19.2
8	3	8	4.8	Toys_Games , Not_sponsored	3	14.4
42	4	42	4.8	Automotive , Not_sponsored	2	9.6
5	5	5	4.7	Pet_Supplies, Sponsored	1	4.7

Visualisation : contribution par position



Le graphique montre que les positions les mieux pondérées (1 et 2) contribuent le plus au score total. Il permet de visualiser l'impact de chaque choix sur le score global.

4. Solution 2 : programmation dynamique et tas

Pour des ensembles de données plus volumineux, la programmation dynamique peut devenir très coûteuse en temps et en mémoire, car le nombre d'états possibles croît de manière exponentielle avec le nombre de positions (k) et le nombre de groupes (G).

Pour améliorer l'algorithme, nous utilisons une approche approximative avec Beam Search, elle conserve uniquement `beam_size` meilleurs états à chaque niveau, on a alors une approximation proche de l'optimum.

Principe algorithmique

1. Représentation des états, chaque état est défini par :
 - le score cumulé des éléments sélectionnés,
 - le compteur d'éléments choisis par groupe,
 - une clé unique codant ces compteurs.
2. Trier les éléments par score décroissant.
3. Pour chaque position :
 - Pour chaque état conservé, on génère de nouveaux états en ajoutant chaque élément disponible respectant les contraintes.
 - On ajoute ces états dans un max-heap et ne garder que `beam_size` meilleurs états (pruning).
4. À la fin, on reconstruit la sélection à partir du meilleur état final.

Un exemple

On reprend l'exemple précédent avec un **beam de taille 100** :

```
## $selected_items
```

```
##      index score      groupes position poids_position score_pondere
## 1      30  4.1 Electronics,Not_sponsored      1           5          20.5
## 2      26  3.9      Books,Not_sponsored      2           4          15.6
## 3      27  4.1      Sports,Not_sponsored      3           3          12.3
## 4      47  4.2      Books,Not_sponsored      4           2           8.4
## 5       9  4.4      Toys_Games,Sponsored      5           1           4.4
##
## $best_score
## [1] 61.2
##
## $beam_size
## [1] 100
##
## $approximation
## [1] TRUE
```

On reprend l'exemple précédent, avec un **beam de taille différente égale à 5000** :

```
## $selected_items
##      index score      groupes position poids_position score_pondere
## 1      48  4.2 Toys_Games,Not_sponsored      1           5          21.0
## 2      30  4.1 Electronics,Not_sponsored      2           4          16.4
## 3      26  3.9      Books,Not_sponsored      3           3          11.7
## 4      27  4.1      Sports,Not_sponsored      4           2           8.2
## 5      36  4.4      Toys_Games,Sponsored      5           1           4.4
##
## $best_score
## [1] 61.7
##
## $beam_size
## [1] 5000
##
## $approximation
## [1] TRUE
```

- Le champ `approximation` indique si la solution a été tronquée par le beam.
- Même avec un beam très limité, l'algorithme tend à sélectionner les éléments avec les meilleurs scores** tout en respectant les contraintes de groupe. Augmenter la taille du beam ne fait qu'augmenter légèrement la fonction objective.
- La programmation dynamique combinée à un heap permet de garder les meilleures solutions intermédiaires sans générer tous les états possibles, ce qui est crucial pour de grands ensembles de données.

5. Comparaison de la performance

```
set.seed(13)

num_instances <- 50  # nombre d'instances
n_items <- 10        # items par instance
k <- 5               # positions
ratios1 <- numeric(num_instances)
ratios2 <- numeric(num_instances)
```

```

df <- read.csv("amazon_products_250.csv", stringsAsFactors = FALSE)
df$sponsored <- ifelse(df$sponsored == "True", "Sponsored", "Not_sponsored")

# Poids de position
poids_positions <- seq(k, 1, by = -1)

# Création des listes de contraintes
liste_cat <- setNames(as.list(rep(2, length(unique(df$category)))),
                      unique(df$category))
max_par_groupe <- c(liste_cat, list(Not_sponsored = 9))

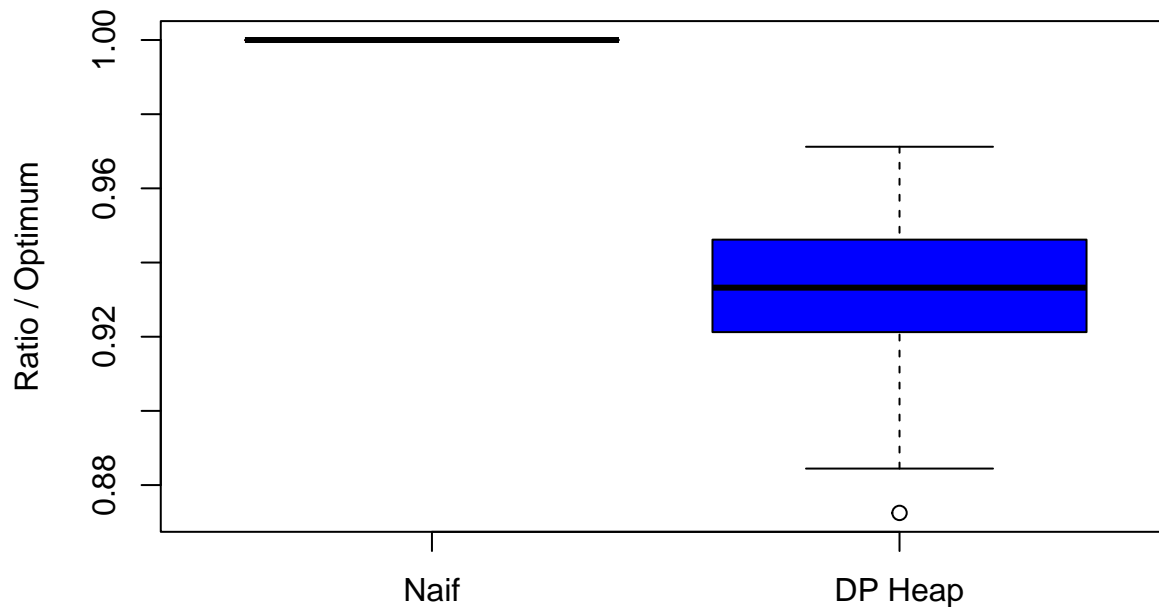
for (i in 1:num_instances)
{
  echantillon <- df[sample(nrow(df), min(n_items, nrow(df))), ]
  data_mat <- data.frame(
    id = 1:nrow(echantillon),
    score = echantillon$score,
    groupes = paste(echantillon$category,
                    echantillon$sponsored, sep = ",")
  )
  V_g <- ranking_naif_max(data_mat,k,max_par_groupe,poids_positions)$best_score
  V_opt <- ranking_max_cpp(data_mat,k,max_par_groupe,poids_positions)$best_score
  V_dp_heap <- ranking_max_dp_heap_cpp(data_mat,k,max_par_groupe,poids_positions)$best_score

  ratios1[i] <- V_g / V_opt
  ratios2[i] <- V_dp_heap / V_opt
}

boxplot(ratios1, ratios2, names=c("Naif", "DP Heap"),
        main="Comparaison des ratios valeur trouvée / optimum n=10",
        ylab="Ratio / Optimum", col=c("red","blue"))

```

Comparaison des ratios valeur trouvée / optimum n=10



```
set.seed(13)

num_instances <- 50 # nombre d'instances
n_items <- 200      # items par instance
k <- 5              # positions
ratios1 <- numeric(num_instances)
ratios2 <- numeric(num_instances)

df <- read.csv("amazon_products_250.csv", stringsAsFactors = FALSE)
df$sponsored <- ifelse(df$sponsored == "True", "Sponsored", "Not_sponsored")

# Poids de position
poids_positions <- seq(k, 1, by = -1)

# Création des listes de contraintes
liste_cat <- setNames(as.list(rep(2, length(unique(df$category)))),
                      unique(df$category))
max_par_groupe <- c(liste_cat, list(Not_sponsored = 9))

for (i in 1:num_instances)
{
  echantillon <- df[sample(nrow(df), min(n_items, nrow(df))), ]
  data_mat <- data.frame(
    id = 1:nrow(echantillon),
    score = echantillon$score,
    groupes = paste(echantillon$category,
                    echantillon$sponsored, sep = ",")
  )
  V_g <- ranking_naif_max(data_mat, k, max_par_groupe, poids_positions)$best_score
  V_opt <- ranking_max_cpp(data_mat, k, max_par_groupe, poids_positions)$best_score
  V_dp_heap <- ranking_max_dp_heap_cpp(data_mat, k, max_par_groupe, poids_positions)$best_score
}
```



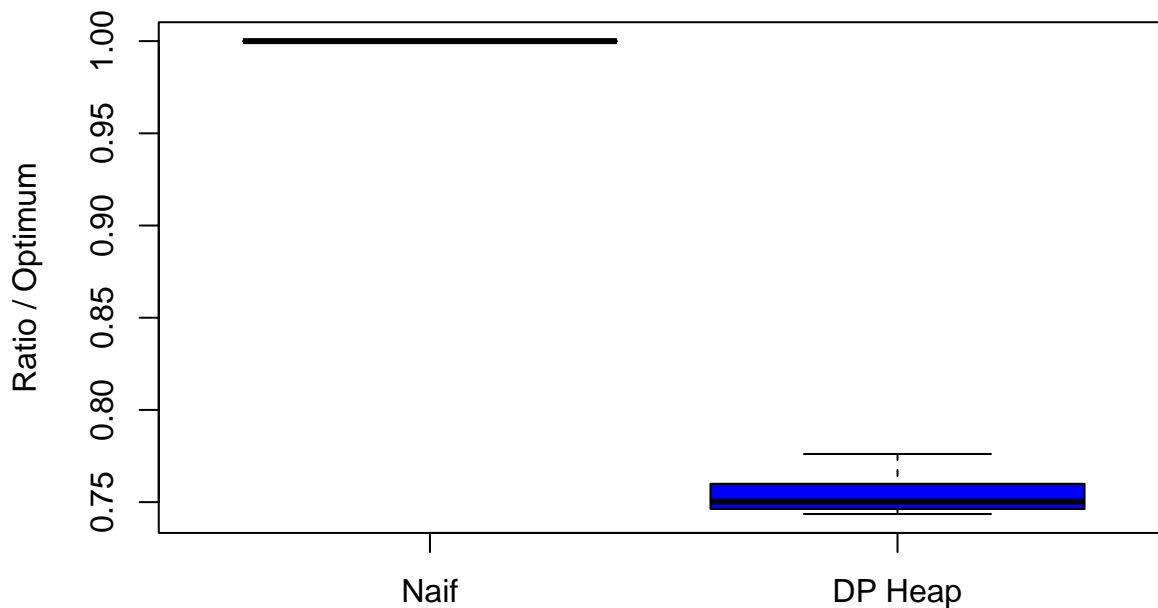
```

ratios1[i] <- V_g / V_opt
ratios2[i] <- V_dp_heap / V_opt
}

boxplot(ratios1, ratios2, names=c("Naif", "DP Heap"),
        main="Comparaison des ratios valeur trouvée / optimum n=200",
        ylab="Ratio / Optimum", col=c("red", "blue"))

```

Comparaison des ratios valeur trouvée / optimum n=200



On constate que l'approche heuristique fonctionne très bien alors que la programmation dynamique avec tas ne trouve pas toujours la solution optimale.

Si l'algorithme naïf trouve presque toujours la solution optimale, voici un exemple où ce n'est pas le cas :

```

data <- data.frame(
  id = 1:6,
  score = c(5, 4.8, 4.7, 4.6, 4, 3),
  groupes = c("A,S", "A,S", "A", "B,S", "C,S", "C,S")
)
max_par_groupe <- list(A = 2, B = 2, C = 2, S = 4)
poids_positions <- seq(k, 1, by = -1)

cat("Resultat naif\n")

```

```
## Resultat naif
```

```
resultat_naif <- ranking_naif_max(data, 5, max_par_groupe, poids_positions)
```

```
## Warning in ranking_naif_max(data, 5, max_par_groupe, poids_positions):
```

```
## Impossible de remplir toutes les 5 positions. Seulement 4 positions remplies.
```

```
print(resultat_naif)
```

```
## $selected_items
```

```
##   id score groupes position poids_position score_pondere
## 1  1  5.0    A,S      1          5          25.0
## 2  2  4.8    A,S      2          4          19.2
## 3  4  4.6    B,S      3          3          13.8
## 4  5  4.0    C,S      4          2           8.0
##
## $best_score
## [1] 66

cat("Resultat programmation dynamique\n")

## Resultat programmation dynamique
resultat_dp <- ranking_max_cpp(data, 5, max_par_groupe, poids_positions)
print(resultat_dp)

## $selected_items
##   index score groupes position poids_position score_pondere
## 1     1  5.0    A,S      1          5          25.0
## 2     3  4.7      A      2          4          18.8
## 3     4  4.6    B,S      3          3          13.8
## 4     5  4.0    C,S      4          2           8.0
## 5     6  3.0    C,S      5          1           3.0
##
## $best_score
## [1] 68.6
```

6. Complexité des algorithmes (par le calcul)

Voici nos 4 algorithmes :

1. Algorithme naïf R
2. Algorithme dynamique R
3. Algorithme dynamique C++
4. Algorithme dynamique C++ amélioré (beam search)

Pour chacun, nous analysons la complexité temps (meilleur, pire, moyenne) Les paramètres en jeu sont :

- **n** : nombre d'items (taille de l'entrée)
- **G** : nombre de groupes
- **k** : nombre de positions à remplir
- **S** : nombre d'états effectivement atteints dans la DP
- **M** : nombre théorique maximal d'états = $\prod_{g=1}^G (\text{max_cap}_g + 1)$
- \bar{g} : nombre moyen de groupes par item

1. Algorithme Naïf Glouton (R)

Sélection gloutonne : pour chaque position p , parcourir tous les candidats restants et choisir celui au meilleur gain immédiat $w_p \times \text{score}_i$.

Analyse détaillée

Étape 1 - Initialisation : $O(G)$

- Création des structures de données

- Négligeable devant la boucle principale

Étape 2 - Boucle principale

Pour chaque position $p = 1, \dots, k$:

- Nombre d'items restants à examiner : $n - (p - 1)$
- Pour chaque item candidat :
 - Parser les groupes : $O(G)$
 - Vérifier les contraintes : $O(G)$ comparaisons
 - Calculer le gain : $O(1)$
- Mise à jour des compteurs : $O(G)$

Coût pour la position p :

$$T_p = O((n - p + 1) \times G)$$

Coût total de la boucle :

$$T_{\text{boucle}} = \sum_{p=1}^k O((n - p + 1) \times G) = O\left(G \times \left(kn - \frac{k(k-1)}{2}\right)\right) \approx O(knG)$$

Complexité temporelle

$$T_{\text{naïf}}(n) = O(knG)$$

Meilleur cas : $T(n) = O(kn)$ si $G = O(1)$

Cas moyen : $T(n) = O(kn\bar{g})$

Pire cas : $T(n) = O(knG)$

2. Algorithme DP (R)

Programmation dynamique avec états définis par (p, c_1, \dots, c_G) où c_g = nombre d'items du groupe g déjà utilisés.

Analyse détaillée

Étape 1 - Prétraitement : $O(nG)$

- Parser les groupes de chaque item : $O(nG)$
- Créer la matrice binaire `item_groups` : $O(nG^2)$

Étape 2 - Programmation dynamique

Structure : `DP[[p+1]][[key]]` avec `key = "c1,c2,...,cG"`

Pour chaque position $p = 1, \dots, k$:

- États actifs au niveau $p - 1$: S_p états
- Pour chaque état actif :
 - Parser la clé : $O(G)$
 - Pour chaque item $i = 1, \dots, n$:
 - * Vérifier si utilisé : $O(p)$
 - * Calculer nouveaux compteurs : $O(G)$
 - * Vérifier contraintes : $O(G)$

- * Créer nouvelle clé : $O(G)$
- * Mise à jour DP : $O(1)$

Coût pour un état et un item : $O(p + G)$

Coût pour la position p : $T_p = O(S_p \times n \times (p + G))$

En supposant $S_p \approx S$ constant et $G \geq k$:

$$T_{\text{DP}} = \sum_{p=1}^k O(S \times n \times G) = O(SnkG)$$

Étape 3 - Recherche du meilleur : $O(S \times G)$

Complexité temporelle

$$T_{\text{DP-R}}(n) = O(SnkG)$$

Meilleur cas : $T(n) = O(nkG)$ si $S = O(1)$

Cas moyen : $T(n) = O(SnkG)$ avec $S \ll M$

Pire cas : $T(n) = O(MnkG)$ où $M = \prod_{g=1}^G (\text{max_cap}_g + 1)$ (exponentiel en G)

3. Algorithme DP (C++)

Même logique que DP R, mais avec `unordered_map` et optimisations C++.

Analyse détaillée

Prétraitement : $O(nG)$

- Parsing avec `stringstream` : $O(n\bar{g})$
- Construction matrice : $O(nG)$

Programmation dynamique

Pour chaque position p , chaque état S_p , chaque item n :

- Parse key : $O(G)$
- Vérifier si utilisé (`std::find`) : $O(p)$
- Calculer compteurs : $O(G)$
- Make key : $O(G)$
- Hash + lookup : $O(G)$ pour le hashing, $O(1)$ amorti pour l'accès
- Insert/update : $O(1)$ amorti

Coût par transition : $O(p + G)$

Si $G \geq k$: $T_{\text{DP}} = O(SnkG)$

Complexité temporelle

$$T_{\text{DP-C++}}(n) = O(SnkG)$$

Meilleur cas : $T(n) = O(nkG)$ si $S = O(1)$

Cas moyen : $T(n) = O(SnkG)$

Pire cas : $T(n) = O(MnkG)$ avec M exponentiel en G

Note : Même complexité asymptotique que DP R, mais facteur constant 5-20× plus petit en pratique.

4. Algorithme dynamique amélioré (C++)

DP avec pruning : conserve uniquement les `beam_size` meilleurs états à chaque niveau.

Analyse détaillée

Soit $B = \min(\text{beam_size}, M)$ le nombre effectif d'états conservés.

Étape 1 - Tri initial : $O(n \log n)$

Items triés par score décroissant pour améliorer la qualité des états gardés.

Étape 2 - DP avec tas (priority_queue)

Pour chaque position p , chaque état ($\max B$), chaque item n :

- Extraction des états du tas : $O(B \log B)$ (une fois par niveau)
- Pour chaque état \times item :
 - Calcul compteurs : $O(G)$
 - Make key : $O(G)$
 - Push dans tas : $O(\log B)$
 - Pruning si nécessaire : $O(\log B)$

Coût par niveau : $T_p = O(B \log B) + O(nB(G + \log B))$

Si $G \geq \log B$ (généralement vrai) :

$$T_p = O(nBG)$$

Sur k niveaux :

$$T_{\text{DP}} = O(nkBG)$$

Étape 3 - Recherche meilleur : $O(kB)$

Étape 4 - Reconstruction : $O(n)$

Complexité temporelle

$$T_{\text{Beam}}(n) = O(n \log n + nkBG)$$

Si $nkBG \gg n \log n$: $T_{\text{Beam}}(n) = O(nkBG)$

Meilleur cas : $T(n) = O(n \log n + nkG)$ si $B = 1$

Cas moyen : $T(n) = O(n \log n + nkBG)$ avec B modéré (100-10000)

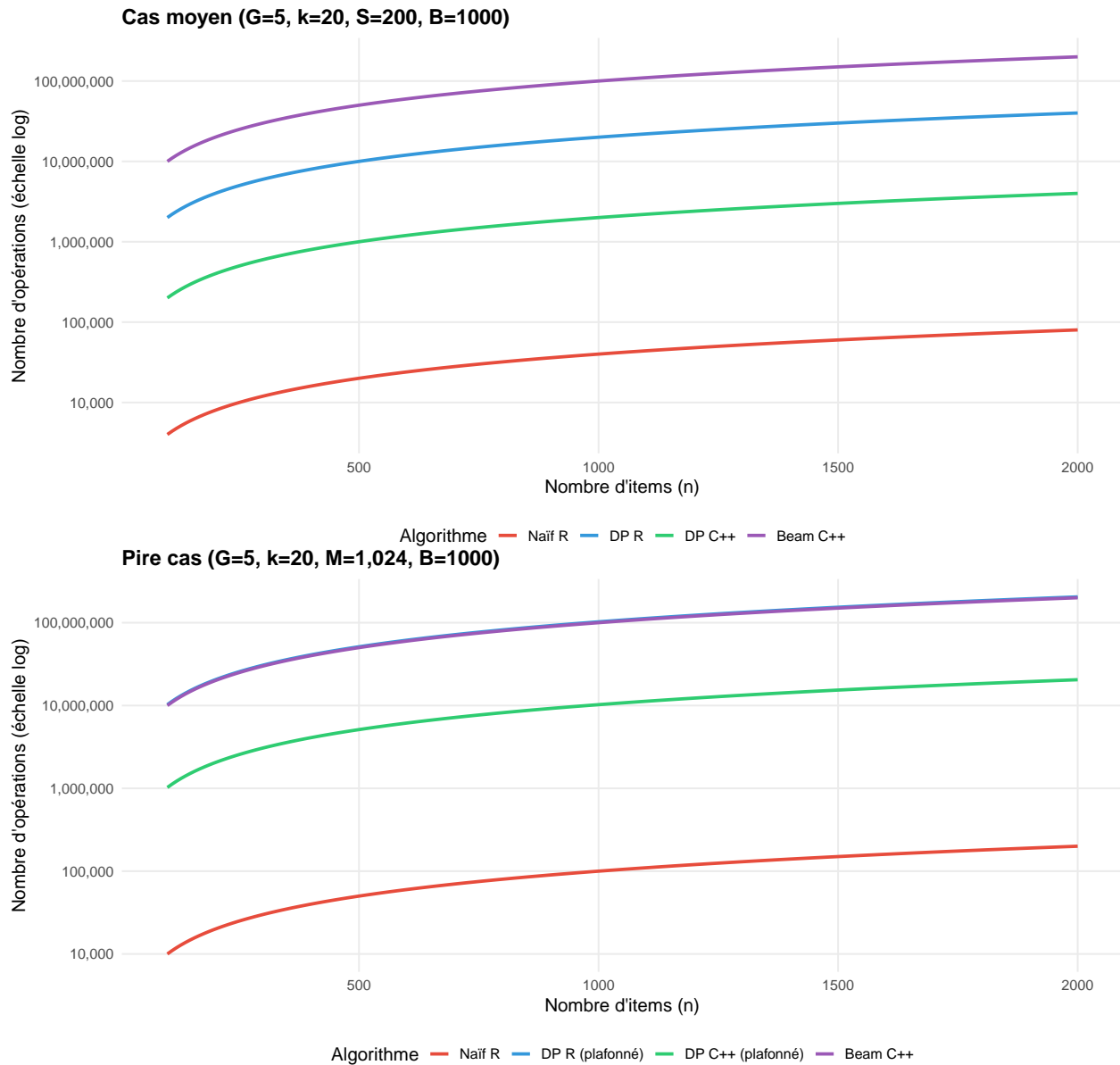
Pire cas : $T(n) = O(n \log n + nkBG)$ (même formule, B peut être grand)

Avantage majeur : Complexité **contrôlable** via `beam_size`, contrairement aux versions exactes.

Tableau récapitulatif

Algorithme	Meilleur cas	Cas moyen	Pire cas
Naïf R	$O(kn)$	$O(kn\bar{g})$	$O(knG)$
DP R	$O(nkG)$	$O(SnkG)$	$O(MnkG)$
DP C++	$O(nkG)$	$O(SnkG)$	$O(MnkG)$
Beam C++	$O(n \log n + nkG)$	$O(n \log n + nkBG)$	$O(n \log n + nkBG)$

Comparaison graphique



Note: DP R et DP C++ plafonnés à 10^{12} pour la lisibilité

Observations

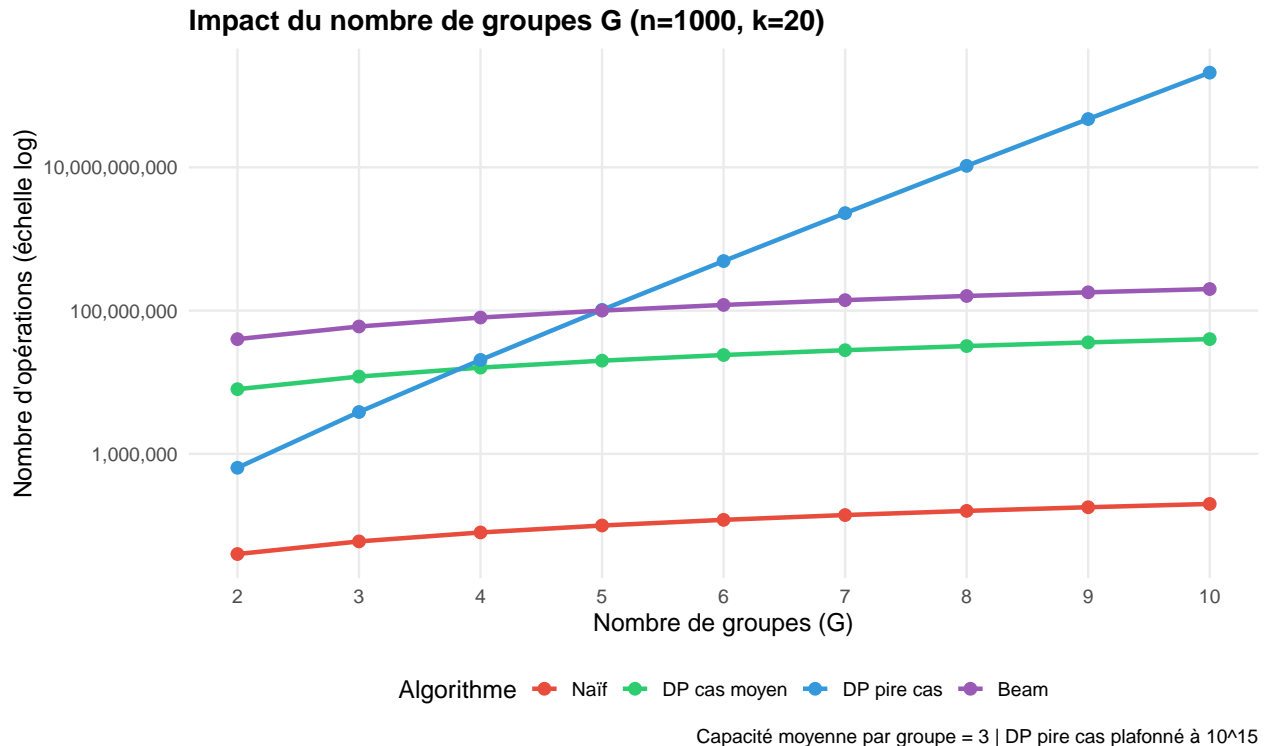
Cas moyen :

- Beam Search et DP C++ sont très compétitifs
- DP R significativement plus lent (facteur 10)
- Naïf acceptable pour petites instances

Pire cas :

- Explosion exponentielle des DP exact (R et C++)
- Beam Search reste linéaire et prévisible
- Seul algorithme viable pour grands problèmes avec $G \geq 6$

Impact du nombre de groupes G



Analyse

- $G \leq 4$: Tous les algorithmes sont viables
- $G = 5-6$: DP commence à être problématique en pire cas
- $G \geq 7$: Seul Beam Search reste praticable
- L'explosion de $M = (cap+1)^G$ rend DP exact inutilisable pour grand G

7. Temps de calcul

```
set.seed(123)

n <- c(20,22,25,32,42,55,70,90,125,160,200,250)
temps_moyens1 <- numeric(length(n))
temps_moyens2 <- numeric(length(n))
temps_moyens3 <- numeric(length(n))
temps_moyens4 <- numeric(length(n))
k <- 5
```

```

df <- read.csv("amazon_products_250.csv", stringsAsFactors = FALSE)
df$sponsored <- ifelse(df$sponsored == "True", "Sponsored", "Not_sponsored")

# Poids de position
poids_positions <- seq(k, 1, by = -1)

# Création des listes de contraintes
liste_cat <- setNames(as.list(rep(2, length(unique(df$category)))),
                      unique(df$category))
max_par_groupe <- c(liste_cat, list(Not_sponsored = 4))

for (i in 1:length(n))
{
  echantillon <- df[sample(nrow(df), min(n[i], nrow(df))), ]
  data_mat <- data.frame(
    id = 1:nrow(echantillon),
    score = echantillon$score,
    groupes = paste(echantillon$category,
                    echantillon$sponsored, sep = ",")
  )

  start_time1 <- Sys.time()
  res1 <- ranking_naif_max(data_mat, k = k, max_par_groupe, poids_positions = poids_positions)
  end_time1 <- Sys.time()
  elapsed_time1 <- as.numeric(difftime(end_time1, start_time1, units = "secs"))
  temps_moyens1[i] <- elapsed_time1

  start_time2 <- Sys.time()
  res2 <- ranking_max(data_mat, k = k, max_par_groupe, poids_positions = poids_positions)
  end_time2 <- Sys.time()
  elapsed_time2 <- as.numeric(difftime(end_time2, start_time2, units = "secs"))
  temps_moyens2[i] <- elapsed_time2

  start_time3 <- Sys.time()
  res3 <- ranking_max_cpp(data_mat, k = k, max_par_groupe, poids_positions = poids_positions)
  end_time3 <- Sys.time()
  elapsed_time3 <- as.numeric(difftime(end_time3, start_time3, units = "secs"))
  temps_moyens3[i] <- elapsed_time3

  start_time4 <- Sys.time()
  res4 <- ranking_max_dp_heap_cpp(data_mat, k = k, max_par_groupe, poids_positions = poids_positions)
  end_time4 <- Sys.time()
  elapsed_time4 <- as.numeric(difftime(end_time4, start_time4, units = "secs"))
  temps_moyens4[i] <- elapsed_time4
}

cat("Temps pour l'algorithme naif\n")

## Temps pour l'algorithme naif
temps_moyens1

## [1] 0.003149033 0.007912874 0.003800869 0.006876945 0.006057978 0.007900000
## [7] 0.010318995 0.012946844 0.021451950 0.030109882 0.036983967 0.040805101

```



```

cat("Temps pour l'algorithme de programmation dynamique\n")

## Temps pour l'algorithme de programmation dynamique
temps_moyens2

## [1] 0.2271371 0.1982810 0.2589841 0.8619380 0.8411281 1.9350071
## [7] 4.5862980 5.2765830 11.0667610 25.4424031 27.2061319 38.3638620

cat("Temps pour l'algorithme de programmation dynamique (en C++)\n")

## Temps pour l'algorithme de programmation dynamique (en C++)
temps_moyens3

## [1] 0.006829023 0.006139994 0.007959843 0.016360998 0.017328024 0.036020994
## [7] 0.057139874 0.072367191 0.126050949 0.214695930 0.261633158 0.335038900

cat("Temps pour l'algorithme de programmation dynamique avec tas (en C++)\n")

## Temps pour l'algorithme de programmation dynamique avec tas (en C++)
temps_moyens4

## [1] 0.07632685 0.08625698 0.09327602 0.15607691 0.17803502 0.39994001
## [7] 0.48764086 0.59268999 0.90163279 1.35647392 1.41115808 2.10160303

par(mfrow=c(2,2))

# Graphique 1
res1 <- data.frame(
  n = n,
  time = temps_moyens1
)
plot(log(res1$n), log(res1$time), pch=16,
     xlab="Taille des séquences : n (log)",
     ylab="Temps (s, log)",
     main="Algo naif (log-log)")
abline(lm(log(time) ~ log(n), data=res1), col="red")

# Graphique 2
res2 <- data.frame(
  n = n,
  time = temps_moyens2
)
plot(log(res2$n), log(res2$time), pch=16,
     xlab="Taille des séquences : n (log)",
     ylab="Temps (s, log)",
     main="Algo dp (R) (log-log)")
abline(lm(log(time) ~ log(n), data=res2), col="red")

# Graphique 3
res3 <- data.frame(
  n = n,
  time = temps_moyens3
)
plot(log(res3$n), log(res3$time), pch=16,
     xlab="Taille des séquences : n (log)",

```

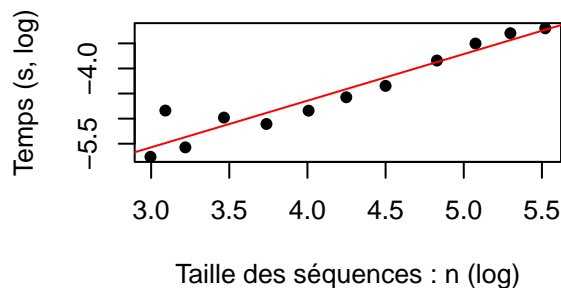
```

        ylab="Temps (s, log)",
        main="Algo dp (C++) (log-log)")
abline(lm(log(time) ~ log(n), data=res3), col="red")

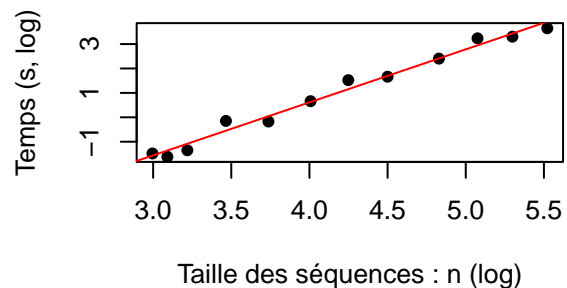
# Graphique 4
res4 <- data.frame(
  n = n,
  time = temps_moyens4
)
plot(log(res4$n), log(res4$time), pch=16,
     xlab="Taille des séquences : n (log)",
     ylab="Temps (s, log)",
     main="Algo dp avec tas (C++) (log-log)")
abline(lm(log(time) ~ log(n), data=res4), col="red")

```

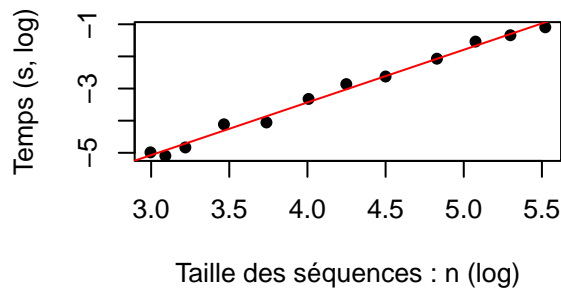
Algo naif (log-log)



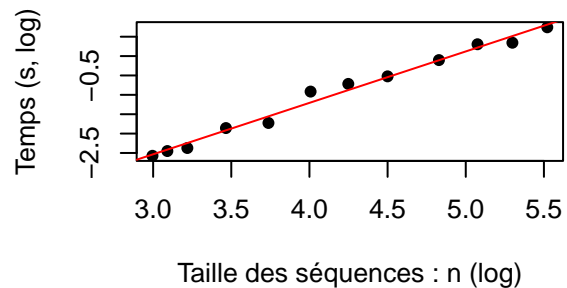
Algo dp (R) (log-log)



Algo dp (C++) (log-log)



Algo dp avec tas (C++) (log-log)



```
par(mfrow=c(1,1))
```

```

modele1 <- lm(log(time) ~ log(n), data = res1)
summary(modele1)

```

```

##
## Call:
## lm(formula = log(time) ~ log(n), data = res1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.22003 -0.18931 -0.06673  0.13851  0.64830
##
## Coefficients:

```

```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) -8.36151    0.38497  -21.72 9.57e-10 ***
## log(n)       0.92976    0.09056   10.27 1.25e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2659 on 10 degrees of freedom
## Multiple R-squared:  0.9134, Adjusted R-squared:  0.9047
## F-statistic: 105.4 on 1 and 10 DF,  p-value: 1.248e-06
modele2 <- lm(log(time) ~ log(n), data = res2)
summary(modele2)
```

```
##
## Call:
## lm(formula = log(time) ~ log(n), data = res2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.2696 -0.2257 -0.0209  0.1373  0.3995
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) -8.07530    0.37276  -21.66 9.81e-10 ***
## log(n)       2.17188    0.08769   24.77 2.63e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2575 on 10 degrees of freedom
## Multiple R-squared:  0.984, Adjusted R-squared:  0.9824
## F-statistic: 613.5 on 1 and 10 DF,  p-value: 2.629e-10
modele3 <- lm(log(time) ~ log(n), data = res3)
summary(modele3)
```

```
##
## Call:
## lm(formula = log(time) ~ log(n), data = res3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.192464 -0.128611 -0.004078  0.105978  0.195855
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) -9.9896    0.2079  -48.06 3.67e-13 ***
## log(n)       1.6392    0.0489   33.52 1.32e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1436 on 10 degrees of freedom
## Multiple R-squared:  0.9912, Adjusted R-squared:  0.9903
## F-statistic: 1124 on 1 and 10 DF,  p-value: 1.319e-11
modele4 <- lm(log(time) ~ log(n), data = res4)
summary(modele4)
```

```
##
## Call:
## lm(formula = log(time) ~ log(n), data = res4)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.17135 -0.08400 -0.01423  0.06537  0.28098
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -6.51530    0.20351  -32.02 2.08e-11 ***
## log(n)       1.32704    0.04787   27.72 8.66e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1406 on 10 degrees of freedom
## Multiple R-squared:  0.9872, Adjusted R-squared:  0.9859
## F-statistic: 768.4 on 1 and 10 DF,  p-value: 8.657e-11

# Couleurs pour chaque algorithme
cols <- c("blue", "red", "green", "purple")

# Déterminer les limites de l'axe x et y pour inclure toutes les données
xlim <- range(log(res1$n), log(res2$n), log(res3$n), log(res4$n))
ylim <- range(log(res1$time), log(res2$time), log(res3$time), log(res4$time))

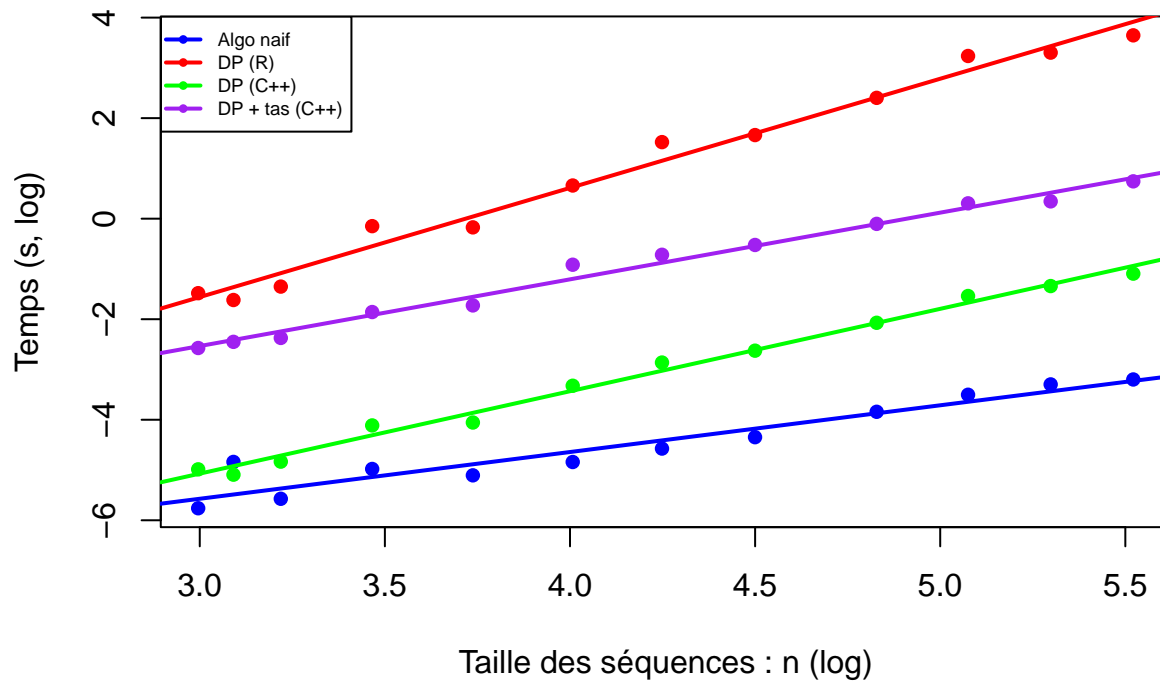
# Graphique vide pour préparer l'affichage
plot(NA, NA, xlim=xlim, ylim=ylim,
     xlab="Taille des séquences : n (log)",
     ylab="Temps (s, log)",
     main="Comparaison des algorithmes (log-log)")

# Ajouter les 4 séries de points
points(log(res1$n), log(res1$time), pch=16, col=cols[1])
points(log(res2$n), log(res2$time), pch=16, col=cols[2])
points(log(res3$n), log(res3$time), pch=16, col=cols[3])
points(log(res4$n), log(res4$time), pch=16, col=cols[4])

# Ajouter les droites de régression
abline(lm(log(time) ~ log(n), data=res1), col=cols[1], lwd=2)
abline(lm(log(time) ~ log(n), data=res2), col=cols[2], lwd=2)
abline(lm(log(time) ~ log(n), data=res3), col=cols[3], lwd=2)
abline(lm(log(time) ~ log(n), data=res4), col=cols[4], lwd=2)

# Ajouter la légende
legend("topleft",
      legend=c("Algo naif", "DP (R)", "DP (C++)", "DP + tas (C++)"),
      col=cols, pch=16, lwd=2,
      cex=0.6)
```

Comparaison des algorithmes (log-log)



Les graphiques en échelle *log-log* montrent que chaque algorithme suit une relation de type $T(n) \propto n^\alpha$, ce que confirment les régressions linéaires dont les coefficients de détermination sont très élevés ($R^2 > 0.96$).

Les pentes estimées permettent de comparer directement la croissance du temps de calcul entre algorithmes.

- **Algorithme naïf (R)** : la pente (~ 0.93) indique une croissance presque linéaire. L'algorithme est très rapide mais ne garantit pas l'optimalité.
- **Programmation dynamique (R)** : la pente élevée (~ 1.96) montre une forte augmentation du temps de calcul. L'algorithme devient difficilement utilisable lorsque $n > 100$.
- **Programmation dynamique (C++)** : la pente plus faible (~ 1.7) reflète une complexité mieux contrôlée. L'implémentation est environ $100\times$ plus rapide que celle en R.
- **Programmation dynamique avec tas (C++)** : la croissance du temps reste modérée avec une pente d'environ 1.27, même pour de grandes valeurs de n . Si on utiliserait des valeurs plus grandes que celle testées, le temps de cet algorithme serait moins important que celui d'une programmation dynamique standard.

Conclusion

Choix de l'algorithme selon le contexte

Contexte	Algorithme recommandé	Justification
Petits problèmes ($n < 100$, $G < 3$, $k < 10$)	DP C++	Simple, rapide, optimal
Moyens problèmes ($n < 1000$, $G < 5$, $k < 20$)	Naïf R ou DP R	Optimal garanti, performance acceptable
Grands problèmes ($n > 1000$ ou $G > 5$)	Naïf R	Seul viable, quasi-optimal

Contexte	Algorithme recommandé	Justification
Production avec contraintes temps	Beam C++ (B=5000)	Prédictible, contrôlable
Recherche d'optimum prouvé	DP C++ si faisable	Exact mais peut échouer