

Projet de Programmation Impérative

K Plus Proches Voisins

Mathevon B., Palisse E.

UJM - L2 Info

Semestre 4

- 1 K fixe / Nombre de points variables
- 2 Nombre de points fixe / K variable
- 3 Conclusion

Mesure du temps de calcul des KPPV en fonction de deux critères :

- Un nombre de K fixe
- Un nombre de points fixe

K prend la valeur suivante : $k = (2^n) + 1$

K allant de 1 a n-1 , max = 10000

N prend les valeurs suivantes :

10-25-50-100-250-500-1000-5000-10000-50000-100000-500000-1000000

Pour chaque valeur de N : 10 fichiers aléatoires générés à partir du script
`data_gen.sh(nb_n, nb_dim, nb_classe)`

Pour chaque valeur de n et chaque valeur de k : calcul du temps $\times 10$

Deux méthodes de tris utilisés pour le tableau :

- Tri à bulles ($O(n^2)$)
- Tri fusion ($O(n \log n)$)

Pour le tri à bulles : $n_{\max} = 10\,000$

Pour le tri fusion : $n_{\max} = 1\,000\,000$

Arbre KD ($O(\log n)$) : $n_{\max} = 1\,000\,000$

Calcul de la moyenne pour chaque méthode de tri et chaque valeur de k/n

Calcul de l'écart-type (affichage sur les courbes)

Scripts utilisés

Temps de calcul :

Usage : script_time.sh

⟨nb_points⟩

⟨valeur_de_k⟩

⟨mode(tab/tab_fusion/arbre)⟩

Moyenne des temps de calcul :

Usage: script_moyenne.sh

⟨type_moy (en fct de k/n)⟩

⟨nb_n si n/nb_k si k⟩

⟨mode⟩

Génération d'un graphique :

Usage: script_graph.sh

⟨en fct de n ou k⟩

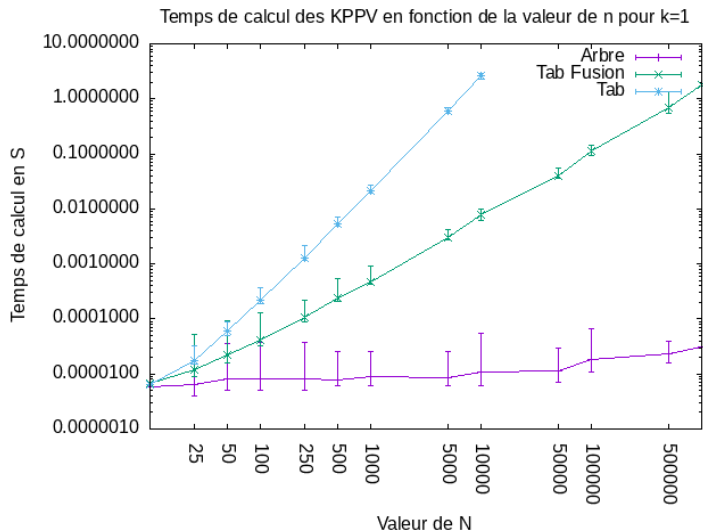
⟨valeur de n/k⟩

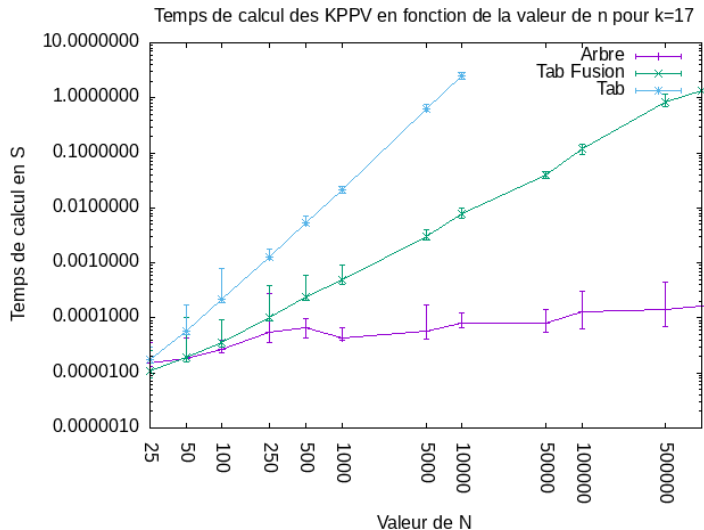
Table des Matières

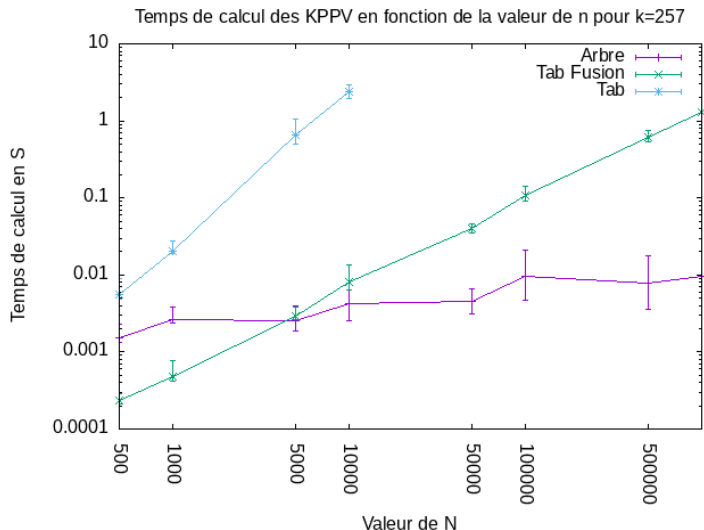
1 K fixe / Nombre de points variables

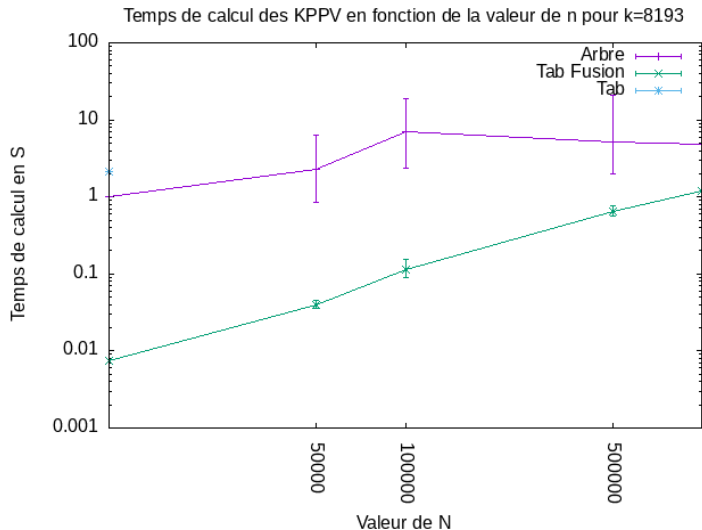
2 Nombre de points fixe / K variable

3 Conclusion









Conclusion K fixe

Tri par tableau fusion : plus le nombre de points (n) est élevé, plus le temps de calcul augmente.

Le temps nécessaire pour trier les données augmente de manière significative lorsque le nombre de points augmente.

Méthode moins efficace pour de grands ensembles de données.

Méthode par arbre kd : le temps de calcul augmente légèrement avec le nombre de points (n).

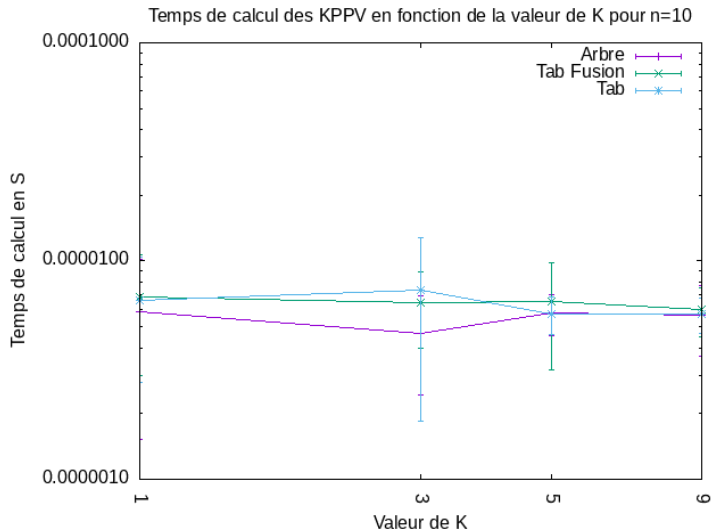
Temps de calcul toujours plus faible pour un $k < 10$. Donc méthode plus efficace pour un k faible peu importe le nombre de points.

Méthode plus efficace pour des ensembles de données plus importants. Mais inversion de la rapidité lorsque le k est très grand.

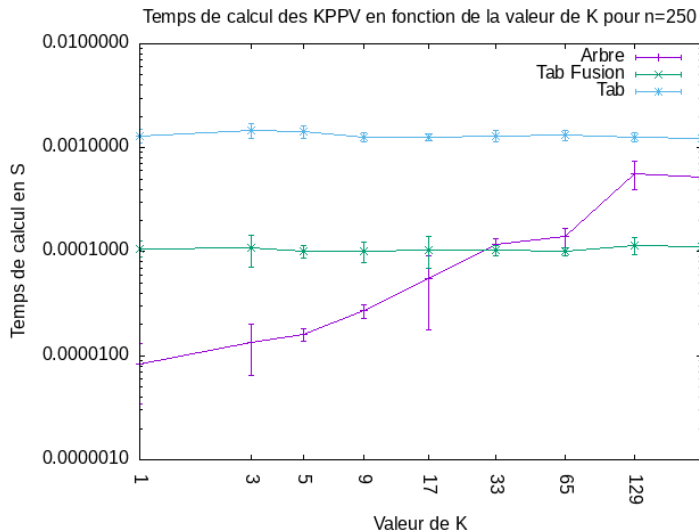
Table des Matières

- 1 K fixe / Nombre de points variables
- 2 Nombre de points fixe / K variable
- 3 Conclusion

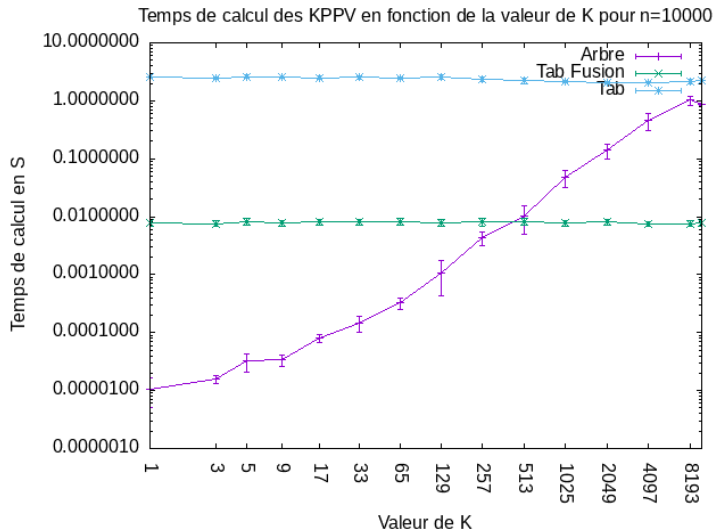
Nombre de points fixe



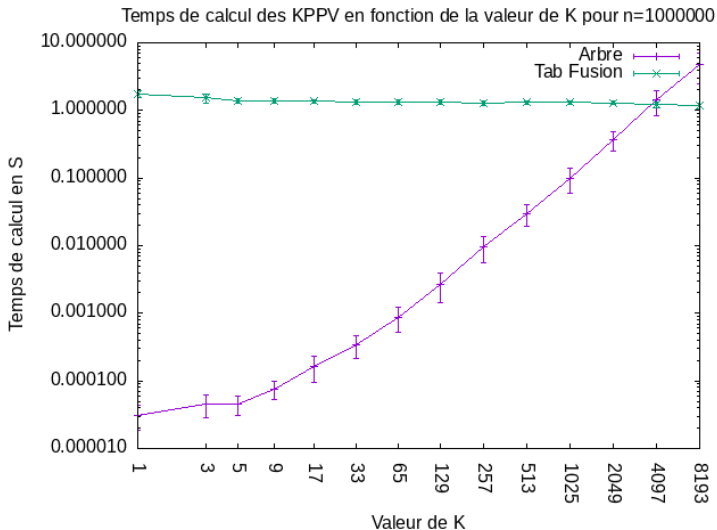
Nombre de points fixe



Nombre de points fixe



Nombre de points fixe



Conclusion pour nombre de points fixe

On voit que le temps de calcul par "tableau" reste constant, peu importe la valeur de K .

La méthode par arbre kd est plus efficace pour un k relativement faible par rapport au nombre de points.

Plus n augmente, plus la valeur de x de l'intersection entre les deux courbes augmente (valeur de k pivot de l'efficacité entre les deux méthodes) car le temps de calcul du tableau augmente.

Table des Matières

- 1 K fixe / Nombre de points variables
- 2 Nombre de points fixe / K variable
- 3 Conclusion

Conclusion

- Pour un k relativement faible par rapport au nombre de points, la méthode par "arbre kd" est plus efficace que la méthode par "tableau".
- Grande valeur de n , faible valeur de k : arbre kd.
- Grande valeur de n , grande valeur de k : tableau.
- Faible valeur de n , faible valeur de k : arbre kd.
- Faible valeur de n , grande valeur de k : tableau.
- On voit également sur les courbes que l'écart-type est plus élevé pour la méthode "arbre kd", car fin de parcours de l'arbre plus aléatoire en fonction du placement des points. Donc temps de calcul plus variable que le tableau.