

# Projet Programmation Impérative

## Calcul efficace du PageRank

Julien Blanchon (AB-15)

16-01-2020

### Résumé

Ceci est un rapport de présentation et de discussion sur une implémentation d'un algorithme de PageRank. Nous allons premièrement un peu de théorie sur l'algorithme PageRank puis ensuite l'architecture de deux implémentations, l'une avec une matrice dite "pleine" (`pagerank_t`) et l'autre avec une matrice dite "creuse" (`pagerank_c`). Les principaux choix et algorithme utilisé seront présentés ainsi que des benchmarks des deux implémentations sur des réseaux tests.

## Table des matières

<b>Introduction:</b>	<b>2</b>
<b>Choix réalisés</b>	<b>2</b>
<b>Raffinages:</b>	<b>3</b>
Consigne : . . . . .	3
Raffinage 0 : Programme Principale . . . . .	4
Raffinage 1 : Vérifier l'intégrité des arguments et les chargés en mémoire . . .	4
Raffinage 2 : Initialiser les paramètres avec les valeurs par défauts . . . . .	5
Raffinage 2 : Mettre à jour les paramètres . . . . .	5
<b>Principaux Algorithme:</b>	<b>5</b>
Calcul et stockage de la matrice $G$ ( <code>Matrice_Pleine</code> ): . . . . .	5
Calcul et stockage de la matrice $G$ ( <code>Matrice_Creuse</code> ): . . . . .	6
Itération Puissance Itérée ( <code>Matrice_Pleine</code> ): . . . . .	6
Itération Puissance Itérée ( <code>Matrice_Creuse</code> ): . . . . .	7
Trie Rapide ( <code>QuickSort</code> ) du vecteur poids: . . . . .	7
<b>Difficulté rencontré, solution et conclusion</b>	<b>7</b>

## Introduction:

L'algorithme de PageRank est un algorithme d'analyse de graphes orientés dont l'objectif est de trier les noeuds selon leurs réputations c'est-à-dire selon leurs popularités (nombre de liens vers le noeud) mais aussi leurs respectabilités (réputations relatives des noeuds liants) (*voir Figure 1*).

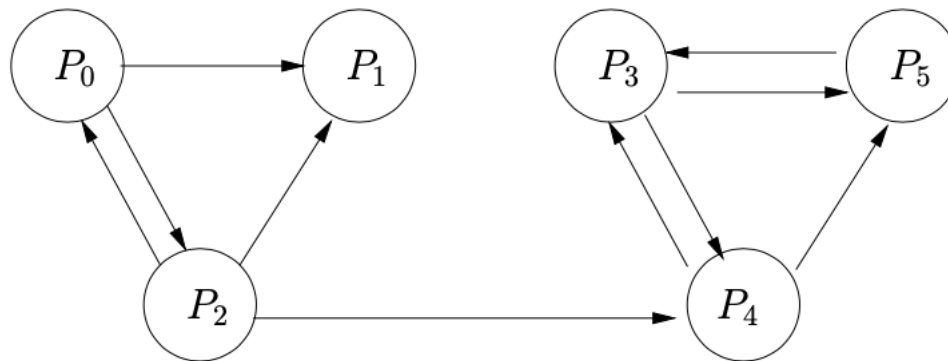


FIGURE 1 – Exemple de graph orienté

Il est au coeur du robot d'indexation de Google **GoogleBot** qui applique ce principe aux graphes des **hyperliens** des pages **html**.

Mathématiquement les réputations des noeuds d'un graph orienté appelé pagerank ou vecteur de poids ( $\pi$ ) peut être approchée (en point fixe) par un algorithme de puissance itérée (*voir Équation 1* calcul de rayon spectral, plus grande valeur propre  $\lambda_{max}$ ) dont l'itération courante est :

$$\pi_{k+1} = \pi_k G \quad (1)$$

Avec  $G$  une matrice dérivée de la matrice adjacente du graph orientée, dont on a normalisé les valeurs et ajouté un coefficient (dumping facteur) (pour ajouter de l'inertie peut être ?). Pour espérer une convergence de l'algorithme il est bien évidemment nécessaire de prendre  $\pi$  non nul !

## Choix réalisés

- Gérer la gestion des arguments avec les exceptions, pour éviter les redondances de code.
- Effectuer l'ensemble des opérations à l'aide de procédures plutôt que de fonction pour éviter de devoir produire des copies des données. Ainsi on utilise une seule et unique matrice  $G$  qui se voit être mise à jour au cours de l'exécution ( $A$  puis  $S$  puis  $G$  ou simplement  $A$ ).
- Utilisation de pointeur pour **Matrice\_Pleine** pour stocker les données dans la **heap** et ainsi pouvoir traiter des réseaux plus grands sans avoir à changer **ulimit** (usage de **Storage Pool** d'Ada).

- On effectue tous les calculs, ainsi que les interactions avec les fichiers dans les modules matrices pour optimiser au mieux les calculs.
- Stockage de certaines constantes pour ne pas avoir à les re-calculer/re-converter (Zero\_D, Un\_D, Taille\_D...).
- Les multiplications matrice-vecteur sont adaptés à la structure du problème. Pour des contraintes de calcul il vaut mieux parcourir les matrices par lignes donc la formule du produit est adaptée de la sorte.
- En **Matrice\_Creuseon** ne stocke pas vraiment la matrice G mais seulement une représentation de celle-ci. C'est-à-dire la matrice Adjacente *A* (**A**) ainsi que le tableau **TabCount** avec le nombre d'éléments non nuls de *A* sur chaque ligne. Ainsi il rapide de déterminer :
  - Les lignes nulles (ainsi que la valeur des coefficients associés qui est  $\frac{1}{Taille} + cst$  stocké sous forme de constante du module).
  - Les coefficients nul (qui sont ceux absents de *A*) dont la valeur vaut une constante *cst*.
  - Les coefficients non nuls dont la valeur vaut  $\frac{1}{Count} + cst$  avec  $\frac{1}{Count}$  stocké directement dans **TabCount**.
- On effectue un Rrie Rapide (**QuickSort**) qui à l'inconvénient d'être récursif mais cependant la profondeur d'appel effective est trop faible pour induire un débordement de la pile.
- De manière générale sauf pour le Trie Rapide on n'utilise pas de fonction récursive. On préfère les fonctions itératives qui évite tout débordement de pile.
- Les programmes de tests ont été implémenter puis abandonner faute de ne plus être adapté à l'interface des modules (et fautes de temps pour les réadapter).

## Raffinages:

### Consigne :

1. L'exécutable **./pagerank** doit comporter 4 argument dont 3 optionnels:
  - **fichier.net** est le chemin vers le fichier décrivant le réseau.
  - **-P** optionnel: Permet d'utiliser le l'implantation **Google\_Naive**. Par défaut (sans **-P**), on lance l'implantation **Google\_Creuse**.
  - **-I <int>** optionnel: Permet de spécifier le nombre maximal d'itérations. Par défaut 150.
  - **-A <float>** optionnel: Permet de spécifier le valeur d' $\alpha$ . Par défaut 0.85.
  - On adopte un programmation défensive pour l'appel de **./pagerank**.
  - **./pagerank -P -I 150 -A 0.90 exemple\_sujet.net**
2. Les résultat du programme seront écrits dans des fichier d'extensions **.p** (pour le fichier pagerank) et **.ord** (pour le fichier poid) avec le même préfixe que le fichier **.net**.
  - **\*.p** Liste le poids des nœuds par ordre décroissant.
  - **\*.ord** Liste l'identifiant des nœuds par poids décroissant (PageRank croissant).

## Raffinage 0 : Programme Principale

---

**Raffinage 0 : Programme Principale**

---

**Input** : stdin("-P -I 150 -A 0.90 exemple\_sujet.net")

**Output** : exemple\_sujet.p:Fichier exemple\_sujet.ord:Fichier

```
1 Vérifier l'intégrité des arguments et les chargés en mémoire
2 if Integre then
3   Calculer la matrice de Google  $G$ 
4   Calculer le vecteur de poids  $\pi$ 
5   Trier  $\pi$  et déterminer le PageRank  $Pk$ 
6   écrire les sorties
7 else
8   Afficher la documentation
```

---

## Raffinage 1 : Vérifier l'intégrité des arguments et les chargés en mémoire

---

**Raffinage 1 : Vérifier l'intégrité des arguments et les chargés en mémoire**

---

**Input** : stdin("-P -I 150 -A 0.90 exemple\_sujet.net")

**Output** : type\_matrice:Integer max\_iter:Integer  $\alpha$ :Float integre:Boolean

```
1 Initialiser les paramètres avec les valeurs par défauts
2 begin
3   Mettre à jour les paramètres
4   exception
5     Max_Iter_Arg_Exception =>
6     | integre := False
7     | Afficher("-I max_iter")
8     Alpha_Arg_Exception =>
9     | integre := False
10    | Afficher("-A alpha")
11    Net_Arg_Exception =>
12    | integre := False
13    | Afficher("chemin/fichier.net")
```

---

## Raffinage 2 : Initialiser les paramètres avec les valeurs par défauts

---

**Raffinage 2** : Initialiser les paramètres avec les valeurs par défauts

---

**Input** :

**Output** : type\_matrice max\_iter  $\alpha$  Integre

```
1 type_matrice := 0
2 max_iter := 150
3  $\alpha$  := 0.85
4 integre := True
```

---

## Raffinage 2 : Mettre à jour les paramètres

---

**Raffinage 1** : Mettre à jour les paramètres

---

**Input** : Argument:String Argument\_Count:Integer

**Output** : type\_matrice:Integer max\_iter:Integer  $\alpha$ :Float integre:Boolean A

```
1 if not( $1 \leq \text{Argument\_Count} \leq 6$ ) then
2   | Raise Argument_Exception
3 I := 1
4 while I  $\neq$  Argument_Count do
5   if Argument(I) := "-P" then
6     | type_matrice := 1
7   else if Argument(I) = "-I" then
8     | max_iter := Argument(I+1)
9     | I := I+1
10  else if Argument(I) = "-A" then
11    |  $\alpha$  := Argument(I+1)
12    | I := I+1
13  else
14    | Ouvrir le fichier
15    | Construire la matrice adjacente A
16  I:=I+1
```

---

## Principaux Algorithmes:

### Calcul et stockage de la matrice $G$ (Matrice\_Pleine):

```
1  — Procédure qui crée la matrice A(G) (Adjacente) et TabCount
2  — à partir du fichier Fichier_Reseau.
3  procedure CalculA(Fichier_Reseau: in Ada.Text_IO.File_Type; G: in out T_Matrice_Pleine;
4                    TabCount: in out T_Vecteur_Plein) is
5    X, Y: Integer;
6  begin
7    Initialiser(G);
8    Initialiser(TabCount);
9    while not( end_of_File(Fichier_Reseau) ) loop
10      Get(Fichier_Reseau, Y);
11      Get(Fichier_Reseau, X);
12      G(X, Y) := Un_D; — Un_D = T_Double(1.0)
13      TabCount(Y) := TabCount(Y) + Un_D;
14    end loop;
15  end CalculA;
```

```

17  — Procedure qui créer la matrice S à partir de A(G) et TabCount.
18  procedure CalculS(G: in out T_Matrice_Pleine; TabCount: in T_Vecteur_Plein) is
19      UnDivTaille : constant T_Double := Un_D/Taille_D;
20      UnDivCount : T_Double;
21  begin
22      for Y in 0..Taille-1 loop
23          if TabCount(Y) = Zero_D then
24              — Ligne vide (Zero): 1/Taille sur tt les elts
25              for X in 0..Taille-1 loop
26                  G(X, Y) := UnDivTaille;
27              end loop;
28          else
29              — Ligne non vide: 1/Count sur les elts non nul
30              UnDivCount := Un_D/TabCount(Y);
31              for X in 0..Taille-1 loop
32                  if G(X, Y) = Un_D then
33                      G(X, Y) := UnDivCount;
34                  end if;
35              end loop;
36          end if;
37      end loop;
38  end CalculS;
39
40  — Procedure qui creer la matrice G.
41  procedure CalculG(alpha: in T_Double; G: in out T_Matrice_Pleine) is
42      Val : constant T_Double := (Un_D-alpha)/Taille_D;
43  begin
44      for Y in 0..Taille-1 loop
45          for X in 0..Taille-1 loop
46              G(X, Y) := alpha*G(X, Y) + Val;
47          end loop;
48      end loop;
49  end CalculG;

```

## Calcul et stockage de la matrice $G$ (Matrice\_Creuse):

```

1  — Procedure qui créer la matrice A(G) (Adjacente) et TabCount
2  — à partir du fichier Fichier_Reseau.
3  procedure CalculA(Fichier_Reseau: in Ada.Text_IO.File_Type; G: in out T_Matrice_Creuse;
4      TabCount: in out T_Vecteur_Plein) is
5      X, Y : Integer;
6      Vcopy : T_Vecteur_Creux;
7      Yold : Integer;
8  begin
9      Initialiser(G);
10     Initialiser(TabCount);
11     while not end_of_File(Fichier_Reseau) loop
12         Get(Fichier_Reseau, Y);
13         Get(Fichier_Reseau, X);
14         if G(Y) = Null then
15             G(Y) := New T_Cellule '(X, Un_D, Null);
16             Vcopy := G(Y);
17         else
18             Vcopy.All.Suivant := New T_Cellule '(X, Un_D, Null);
19             Vcopy := Vcopy.All.Suivant;
20         end if;
21         TabCount(Y) := TabCount(Y) + Un_D;
22     end loop;
23     Free(Vcopy);
24
25 end CalculA;
26
27 — Procedure qui pre-calcul les  $\alpha*(1/Count)+(1-\alpha)/Taille$ 
28 procedure MiseAJourTabCount(TabCount: in out T_Vecteur_Plein; alpha: in T_Double) is
29     Val : constant T_Double := (Un_D-alpha)/Taille_D;
30 begin
31     for Y in 0..Taille-1 loop
32         if TabCount(Y) /= Zero_D then
33             TabCount(Y) := (alpha/TabCount(Y)) + Val;
34         end if;
35     end loop;
36 end MiseAJourTabCount;

```

## Itération Puissance Itérée (Matrice\_Pleine):

```

1  — Fonction qui effectue une itération.
2  function Iteration(G: in T_Matrice_Pleine; Pi: in T_Vecteur_Plein) return T_Vecteur_Plein is
3      Pinew : T_Vecteur_Plein;
4  begin
5      Initialiser(Pinew);
6      for Y in 0..Taille-1 loop
7          for X in 0..Taille-1 loop
8              Pinew(X) := Pinew(X) + Pi(Y)*G(X, Y);
9          end loop;
10     end loop;
11     return Pinew;
12 end Iteration;

```

## Itération Puissance Itérée (Matrice\_Creuse):

```
1  -- Fonction qui effectue une itération.
2  function Iteration(G: in T_Matrice_Creuse; Pi: in T_Vecteur_Plein; alpha: in T_Double; TabCount: in
   T_Vecteur_Plein) return T_Vecteur_Plein is
3      Pinew : T_Vecteur_Plein;
4      Curseur : T_Vecteur_Creux;
5      Val : constant T_Double := (Un_D-alpha)/Taille_D;
6      ValVide : constant T_Double := alpha/Taille_D + Val;
7      TabCountVal : T_Double;
8  begin
9      Initialiser(Pinew);
10     for Y in 0..Taille-1 loop
11         Curseur := G(Y);
12         TabCountVal := TabCount(Y);
13         if TabCountVal = Zero_D then
14             for X in 0..Taille-1 loop
15                 Pinew(X) := Pinew(X) + Pi(Y)*ValVide;
16             end loop;
17         else
18             for X in 0..Taille-1 loop
19                 if Curseur /= Null and then Curseur.All.Indice = X then
20                     Pinew(X) := Pinew(X) + Pi(Y)*TabCountVal;
21                     Curseur := Curseur.All.Suivant;
22                 else
23                     Pinew(X) := Pinew(X) + Pi(Y)*Val;
24                 end if;
25             end loop;
26         end if;
27         Curseur := Null;
28         Free(Curseur);
29     end loop;
30     return Pinew;
31 end Iteration;
```

## Trie Rapide (QuickSort) du vecteur poids:

```
1  -- Procedure qui trie Pi et qui garde dans Index des permutations du trie
2  procedure TrierPi(Pi: in out T_Vecteur_Plein; Index: in out T_Vecteur_PleinInteger) is
3  begin
4      for K in 0..Taille-1 loop
5          Index(K) := K;
6      end loop;
7      Tri_Rapide(Pi, Index, 0, Taille-1);
8  end TrierPi;
9
10 -- Tri relatif des elts du tableau T entre premier et dernier avec le pivot.
11 procedure Partitionner(T: in out T_Vecteur_Plein; Index: in out T_Vecteur_PleinInteger;
12     premier: in Integer; dernier: in Integer; pivot: in out Integer) is
13 begin
14     J: Integer;
15     Echanger(T, pivot, dernier);
16     Echanger(Index, pivot, dernier);
17     J := premier;
18     for I in premier..dernier-1 loop
19         if T(I) >= T(dernier) then -- >= pour un trie décroissant.
20             Echanger(T, I, J);
21             Echanger(Index, I, J);
22             J := J + 1;
23         end if;
24     end loop;
25     Echanger(T, dernier, J);
26     Echanger(Index, dernier, J);
27     pivot := J;
28 end Partitionner;
29
30 -- Tri globale du tableau T entre premier en dernier
31 procedure Tri_Rapide(T: in out T_Vecteur_Plein; Index: in out T_Vecteur_PleinInteger;
32     premier: Integer; dernier: Integer) is
33 begin
34     pivot : Integer;
35     if premier < dernier then
36         pivot := premier;
37         Partitionner(T, Index, premier, dernier, pivot);
38         Tri_Rapide(T, Index, premier, pivot-1);
39         Tri_Rapide(T, Index, pivot+1, dernier);
40     end if;
41 end Tri_Rapide;
```

## Difficulté rencontré, solution et conclusion

La principale difficulté a été d'optimiser le programme que ce soit en temps ainsi qu'en mémoire. Pour cela on a préféré le plus souvent calculer plutôt que stocker car certains

réseaux (en particulier **Linux**) sont très gros et pas forcément très dense. Et ils dépasseraient la mémoire accessible. La solution a donc été de stocker une représentation de la matrice  $G$  à savoir la matrice adjacente **A** ainsi que le tableau **TabCount**. Il reste encore beaucoup de travail dans l'optimisation de l'algorithme et de son implémentation.

Ensuite il a aussi été très difficile de bien structurer les différentes procédure et fonction dans des modules (/package) cohérents. Faute de temps ce travail n'est pas correctement abouti.

En conclusion il reste encore beaucoup de choses à faire malheureusement je n'ai pas assez de temps pour terminé comme j'aimerais ce projet. Il semblerait que l'algorithme de puissance itérée ne soit vraiment pas bien adaptée aux trop gros réseaux (tel que **Linux**) il faudrait réfléchir à une autre méthode pour approcher le vecteur poids.