

Présentation du projet

Wikipédia:

- Encyclopédie gratuite en ligne
- Une encyclopédie par langue, avec ses propres règles
- Chaque article peut en citer d'autres



Jeu de Wikipédia:

Passer d'une page à une autre en utilisant uniquement les liens bleus

Exemples d'articles





Estatistiques

Page d'estatistiques		
Noumbe d'articles	5 031	
	3 031	
Pages (Toutes les pages du wiki, y compris les pages de discussion, les redirections, etc.)	10 638	
Fichiers téléversés	0	
Estatistiques d's éditiouns		
Modifications de pages depuis l'installation de Wikipedia	220 731	
Nombre moyen de modifications par page	20,75	
Estatistique des féchouneus		
Utilisateurs inscrits (liste des bénoums)	14 229	
Utilisateurs actifs (liste des bénoums) (Utilisateurs ayant fait au moins une action durant les 30 derniers jours)	17	
Robots (liste des bénoums)	23	
Counétablles d'ouovrage (liste des bénoums)	2	
Administrateurs d'interface (liste des bénoums)	0	
Bureaucrates (liste des bénoums)	0	
Masqueurs de modifications (liste des bénoums)	0	
Stewards (liste des bénoums)	0	
Créateurs de comptes (liste des bénoums)	0	
Importateurs (liste des bénoums)	0	
Importateurs transwiki (liste des bénoums)	0	
Exemptés de blocage IP (liste des bénoums)	0	
Vérificateurs d'utilisateurs (liste des bénoums)	0	
Utilisateurs bloqués de l'outil d'informations sur une IP (liste des bénoums)	0	
Utilisateurs confirmés (liste des bénoums)	0	
Autres statistiques		
Mots dans toutes les pages de contenu	429 394	

Objectifs



Crawler un Wikipédia complet



Trouver un chemin optimal entre toutes les pages Wikipédia



Trouver les cas impossibles



Interpréter les résultats pour définir des stratégies



Jouer avec les poids pour se rapprocher des vraies conditions

Méthodologie et outils

Python 3.10

• Crawler "homemade" en Python

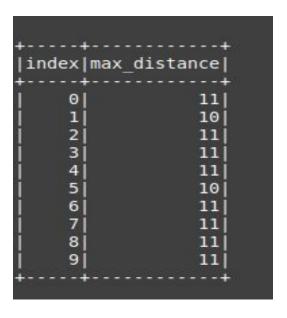
Pyspark(SQL)

Pickle pour stocker les données (pourra être remplacé)

• Matrice d'adjacence pour le graphe

Résultats

Distances depuis chaque article



```
Distances from Oûtriche
|dist|count|
         855
       2558
       1425
Distances from 1322
|dist|count|
         855
       1612
         354
   7 |
8 |
9 |
10 |
|max(max distance)|
```

Résultats

Nombre de pages atteignables

+	++	
Name Reachable		
++		
Câtel_de_Vinchennes		
Edmund_Blampied	4885	
1416	4885	
Rimeus	4885	
1382	4885	
Maûsolée	4885	
1656	4885	
Bailliage_dé_Gùer	4885	
24_Janvyi	4885	
1952	4885	
Jèrri	4885	
27_Novembre	4885	
Esope	4887	
26_Août	4885	
Ph'lippe_Duron	4885	
Berlanga_de_Duero	4886	
25_Juilet	4885	
777	4885	
Nouormandie	4885	
Quien	4885	
++		

Résultats

Chemins les plus courts entre deux articles

```
5 ['Crouésade', 'Jèrriais', 'Langue', 'Laungue construite', 'Angllais']
No path from Clloque to Abejar
 ['757', 'IXe s.', '854']
 ['Ouothelle (frouque)', 'Page dé garde', 'Giographie', 'Liste des pays du monde', 'Kosovo']
  ['Afrique-du-Su', 'Monnaie', 'Liste des dus de Normaundie', '927']
 ['Tokyo', 'Japon', 'Liste des pays du monde', 'Irlaunde']
  ['1890', 'XIXe s.', 'Janvvi', '6 Janvvi']
  ['Annaées 1360', 'Page dé garde', 'Fraunce', 'Crna Gora']
 ['Cannot', 'Page dé garde', 'Mouogeâle', 'Viande']
  ['1645', 'Normaundie', 'Semanne', 'Djödi']
  ['Trochet', 'Page dé garde', 'Aunglléterre', '1603', 'Annaées 1620']
  ['Cité', 'Ville', 'Liste des capitales', 'Muscade', 'Êpice', 'Moutarde']
  ['Histouère', 'Page dé garde', 'Baêtes', 'Mammifère']
  ['Langue romanne', 'Laungue', 'Français']
  ['Mongolīe', 'Page dé garde', 'Mouogeâle', 'Pain', 'Sexualité', 'Contraception', 'Capote']
  ['Scooby-Doo', 'Améthique', 'Ûrope', 'Mathe']
 ['1623', 'Page dé garde', 'Islam', 'Jésû']
  ['Capote', 'Page dé garde', 'Islam', 'Jésû', 'Fei', 'Dieu']
  ['1905', 'Page de garde', 'Afrique', 'Vert Cap']
  ['Sŷmbole', 'Ûro', 'Espangne']
```

Défis rencontrés



Pas de dataset public de taille raisonnable



Le crawling est très long et incertain

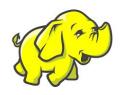


Messages d'erreurs de Java

Défis rencontrés



Pistes d'amélioration



Utiliser HDFS pour le stockage distribué



Remplacer pickle par Parquet



Utiliser un Wikipédia plus grand