



الكلية متعددة التخصصات الناحور

ⵜⴰⵎⵓⵔⵜ ⵜⴰⵎⵓⵔⵜ ⵜⴰⵎⵓⵔⵜ ⵜⴰⵎⵓⵔⵜ | ⵎⵓⵔⵓⵔ

Faculté Pluridisciplinaire de Nador

MASTER "SCIENCES DES DONNÉES ET SYSTÈMES.
INTELLIGENTS (SDSI)"

Analyse de logs avec Pig

Réalisé par :

Meryem Ouafi

Wiam CHOUKOUD

Youssra Andoulla

Encadré par :

Pr.Mourad fariss

ANNÉE UNIVERSITAIRE 2023 - 2024

Table des matières

1	Introduction	2
2	Objectif du projet	2
3	Les technologies utilisees	3
3.1	Hadoop	3
3.2	HDFS (Hadoop Distributed File System)	3
3.3	Apache Pig	3
4	Les donnees qui l'on utilisent :	4
5	Partie pratique :	4
5.1	Importation des donnes dans hdfs :	4
5.2	Lancer pig	5
5.3	Faire des analyses :	5
5.3.1	analyse 1 :Compter le nombre de requêtes pour chaque adresse IP	5
5.3.2	Analyse 2 : de la répartition des codes de statut HTTP	6
5.3.3	Analyse3 :Analyse des navigateurs utilisés :	6
5.3.4	Analyse4 : Trouver le code d'état le plus fréquente :	7
5.3.5	Analyse 5 :Trouver le code d'état le moins fréquente :	7
5.3.6	Analyse 6 : Calculer le nombre total de requêtes par code d'état :	7
5.3.7	Analyse des requêtes par pays	7
6	Conclusion	8

1 Introduction

L'analyse des logs de serveur web est une pratique cruciale dans le monde de l'informatique et du web. Ces fichiers journaux, générés automatiquement par les serveurs web, enregistrent chaque interaction entre le serveur et les utilisateurs. Ils contiennent une mine d'informations sur les requêtes HTTP, les adresses IP des visiteurs, les navigateurs utilisés, les erreurs rencontrées et bien plus encore. L'objectif de cette analyse est de comprendre le comportement des utilisateurs, de diagnostiquer les problèmes de performance, d'identifier les menaces potentielles à la sécurité et d'optimiser l'expérience utilisateur. En examinant les logs de serveur web, les administrateurs système et les spécialistes du marketing peuvent obtenir des informations précieuses sur la popularité du site, les pages les plus visitées, les sources de trafic, les erreurs fréquentes, et bien d'autres aspects essentiels à la gestion efficace d'un site web.



figure 1 :L'analyse des logs

2 Objectif du projet

L'objectif de ce projet est d'utiliser Apache Pig pour analyser des logs de serveur web simples. Nous chercherons à extraire des informations significatives à partir de ces données brutes afin de comprendre le comportement des utilisateurs, diagnostiquer les problèmes potentiels de performance du serveur, et optimiser l'expérience utilisateur. En explorant les différentes analyses que nous pouvons effectuer, nous visons à démontrer comment Apache Pig peut être un outil puissant pour traiter et analyser les logs de serveur web, offrant ainsi des insights précieux pour la gestion et l'optimisation des sites web.



figure 2 :Les fichiers de logs

3 Les technologies utilisees

3.1 Hadoop

Nous utiliserons la plateforme Hadoop pour le stockage et le traitement distribué des données. Hadoop permet de gérer de grands volumes de données de manière évolutive et offre un environnement robuste pour exécuter des analyses sur des clusters de serveurs.

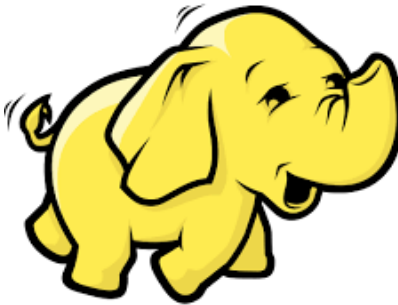


figure 3 : Icone de Hadoop

3.2 HDFS (Hadoop Distributed File System)

Nous stockerons les logs de serveur web dans le système de fichiers distribué Hadoop (HDFS). HDFS permet de stocker les données de manière distribuée sur plusieurs nœuds, offrant ainsi une haute disponibilité et une tolérance aux pannes.



figure 4 :Icone de HDFS

3.3 Apache Pig

Nous utiliserons Apache Pig, un langage de traitement de données parallèle pour Hadoop, pour analyser les logs de serveur web. Pig offre une syntaxe simple et expressive pour effectuer des transformations de données complexes sur de grands ensembles de données de manière efficace.



figure 5 :Icône de Apache Pig

4 Les donnees qui l'on utilisent :

Les données que nous utilisons dans ce projet proviennent des logs d'un serveur web. Chaque ligne de données représente une interaction entre le serveur web et un utilisateur. Voici les différentes colonnes présentes dans nos données :

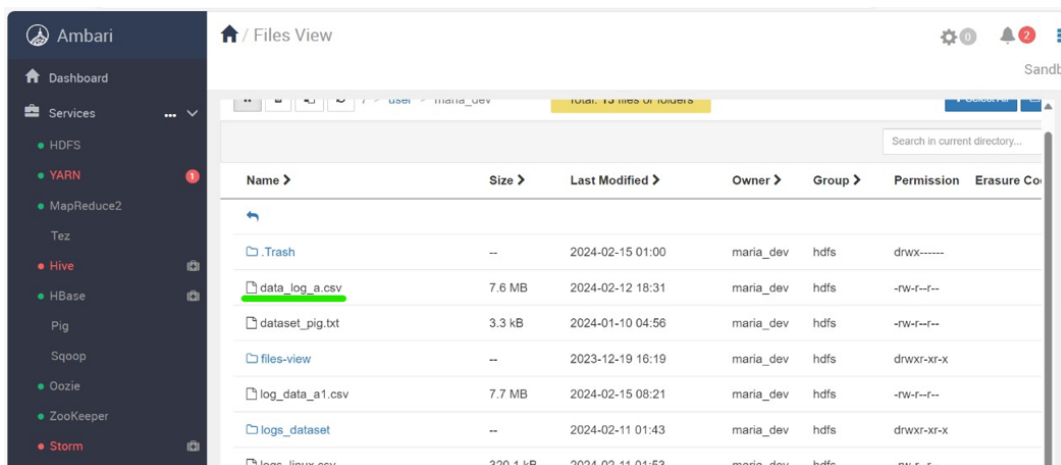
- IP : L'adresse IP de l'utilisateur qui a effectué la requête.
- Timestamp : L'horodatage de la requête, indiquant quand elle a été effectuée.
- Request : Le type de requête HTTP effectuée, par exemple GET, POST, DELETE, etc.
- Navigateur : Le navigateur web utilisé par l'utilisateur pour accéder au site.
- Plateforme : Le système d'exploitation ou la plateforme utilisée par l'utilisateur, par exemple Android, Windows, iOS, etc.
- Country : Le pays à partir duquel la requête a été effectuée.
- Statuicode : Le code de statut HTTP retourné par le serveur en réponse à la requête, indiquant si la requête a réussi, a été redirigée ou a rencontré une erreur.

Ces données fournissent des informations essentielles sur les interactions des utilisateurs avec le site, telles que les caractéristiques de leurs appareils, les pays d'où proviennent les visiteurs, et les éventuelles erreurs rencontrées lors de l'accès au site.

5 Partie pratique :

5.1 Importation des donnees dans hdfs :

Importation de dataset dans hdfs :



5.2 Lancer pig

Lancer pig :

```
(maria_dev@sandbox-hdp ~)$ pig
WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX.
24/02/14 21:29:07 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
24/02/14 21:29:07 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
24/02/14 21:29:07 INFO pig.ExecTypeProvider: Trying ExecType : TEZ_LOCAL
24/02/14 21:29:07 INFO pig.ExecTypeProvider: Trying ExecType : TEZ
24/02/14 21:29:07 INFO pig.ExecTypeProvider: Picked TEZ as the ExecType
2024-02-14 21:29:08,089 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0.3.0.1.0-187 (rUnversioned directory) compiled Sep 19 2018, 10:13:33
2024-02-14 21:29:08,091 [main] INFO org.apache.pig.Main - Logging error messages to: /home/maria_dev/pig_1707946148087.log
2024-02-14 21:29:09,166 [main] INFO org.apache.pig.impl.util.Utils - Default bootup file /home/maria_dev/.pigbootup not found
2024-02-14 21:29:12,817 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://sandbox-hdp1.hortonworks.com:8020
2024-02-14 21:29:25,450 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-default-077551ee-810a-4167-bfd8-b6de9a89d0f2
2024-02-14 21:29:25,451 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
grunt>
```

Charger les donnees de hdfs a pig :

```
grunt> log_a = LOAD 'data_log_a.csv' USING PigStorage(',') AS (IP: chararray, Timestamp: chararray, Request: chararray, Navigateur: chararray, Plateforme: chararray, Country: chararray, Statutcode: int);
grunt> describe log_a;
```

Afficher la description :

```
log_a: (IP: chararray, Timestamp: chararray, Request: chararray, Navigateur: chararray, Plateforme: chararray, Country: chararray, Statutcode: int)
grunt>
```

5.3 Faire des analyses :

5.3.1 analyse 1 : Compter le nombre de requêtes pour chaque adresse IP

Regrouper les logs par adresse IP pour compter le nombre de requêtes pour chaque adresse IP

```
grunt> grouped_logs = GROUP log_a BY IP;
grunt> request_count_per_ip = FOREACH grouped_logs GENERATE group AS IP, COUNT(log_a) AS RequestCount;
grunt> dump request_count_per_ip;
```

```
(233.178.194.204,1)
(233.180.194.183,1)
(233.183.213.102,1)
(233.184.123.172,1)
(233.186.132.221,1)
(233.186.230.183,1)
(233.187.217.143,1)
(233.190.164.134,1)
(233.190.173.225,1)
(233.191.190.125,1)
(233.192.125.101,1)
(233.192.155.138,1)
(233.197.218.248,1)
(233.198.226.187,1)
(233.199.117.233,1)
(233.200.137.100,1)
(233.201.126.234,1)
(233.202.223.135,1)
(233.202.231.210,1)
(233.203.184.251,1)
(233.205.106.178,1)
(233.205.163.163,1)
(233.207.189.233,1)
(233.211.111.183,1)
(233.217.191.213,1)
(233.217.194.215,1)
(233.226.188.222,1)
(233.228.135.121,1)
(233.229.172.141,1)
(233.231.159.119,1)
(233.232.248.106,1)
(233.233.153.149,1)
(233.239.232.230,1)
(233.240.143.205,1)
(233.240.167.128,1)
(233.243.132.182,1)
(233.245.246.140,1)
(233.247.179.213,1)
(233.247.191.213,1)
(233.252.165.215,1)
(233.252.232.142,1)
(234.102.108.209,1)
(234.104.164.203,1)
(234.110.147.244,1)
```

```
(255.196.179.131,1)
(255.197.213.124,1)
(255.197.215.145,1)
(255.198.191.188,1)
(255.199.130.133,1)
(255.199.162.131,1)
(255.201.125.232,1)
(255.203.198.166,1)
(255.204.142.110,1)
(255.204.251.166,1)
(255.208.117.186,1)
(255.209.119.222,1)
(255.211.195.147,1)
(255.213.217.210,1)
(255.213.248.161,1)
(255.214.196.207,1)
(255.215.144.136,1)
(255.217.194.135,1)
(255.217.209.138,1)
(255.218.157.250,1)
(255.222.142.221,1)
(255.224.227.167,1)
(255.225.120.199,1)
(255.225.249.166,1)
(255.225.253.193,1)
(255.226.128.150,1)
(255.226.145.235,1)
(255.226.168.224,1)
(255.231.140.144,1)
(255.231.170.130,1)
(255.233.202.161,1)
(255.234.181.117,1)
(255.235.101.191,1)
(255.239.112.111,1)
(255.239.156.108,1)
(255.244.241.173,1)
(255.249.144.239,1)
(255.251.139.193,1)
(255.252.106.153,1)
(255.254.102.223,1)
(255.254.249.159,1)
(255.255.230.185,1)
(255.255.255.207,1)
```

5.3.2 Analyse 2 : de la répartition des codes de statut HTTP

Examiner la répartition des codes de statut HTTP dans les logs du serveur web. En regroupant les logs par code de statut HTTP

```
(404,42867),2)
grunt> grouped_status_codes = GROUP log_a BY Statuicode;
grunt>
```

```
grunt> status_code_counts = FOREACH grouped_status_codes GENERATE group AS StatusCode, COUNT(log_a) AS Re
grunt>
```

les utilisateurs accèdent principalement au site web à partir d'appareils Android et Windows, avec respectivement 40 162 et 39 810 requêtes. Cependant, on observe également une présence significative d'utilisateurs de Mac OS X, bien que moins importante en termes de volume, avec 9 938 requêtes.

```
(404,42867)
(500,6370)
(501,1)
(502,26007)
(503,24754)
(504,1)
```

5.3.3 Analyse3 :Analyse des navigateurs utilisés :

la répartition des navigateurs utilisés par les utilisateurs pour accéder au site web.

```
grunt> grouped_platforms = GROUP log_a BY Plateforme;
grunt> platform_request_counts = FOREACH grouped_platforms GENERATE group AS Platform, COUNT(log_a) AS RequestCount;
grunt>
```

```
grunt> DUMP platform_request_counts;
```

```
grunt> DUMP status_code_counts;
```

le code de statut HTTP le plus fréquent est 404, avec un total de 42 867 occurrences. Les codes 502 et 503 sont également significativement présents, suggérant potentiellement des erreurs ou des problèmes de connectivité sur le serveur. les utilisateurs accèdent principalement au site web à partir d'appareils Android et Windows, avec respectivement 40 162 et 39 810 requêtes. Cependant, on observe également une présence significative d'utilisateurs de Mac OS X, bien que moins importante en termes de volume, avec 9 938 requêtes. Cette répartition des plateformes utilisées souligne l'importance de l'optimisation de l'expérience utilisateur pour une variété d'appareils et

de systèmes d'exploitation.

```
(Android,40162)
(Windows,39810)
(Mac OS X,9938)
```

5.3.4 Analyse4 : Trouver le code d'état le plus fréquente :

```
grunt> most_frequent_status = FOREACH (GROUP log_a ALL) GENERATE MAX(log_a.Statucode) AS most_frequent_status;
grunt>
grunt>
grunt>
grunt> dump most_frequent_status;
```

```
(504)
grunt>
```

5.3.5 Analyse 5 :Trouver le code d'état le moins fréquente :

```
grunt> least_frequent_status = FOREACH (GROUP log_a ALL) GENERATE MIN(log_a.Statucode) AS least_frequent_status;
grunt>
grunt>
grunt> dump least_frequent_status;
```

```
(404)
grunt>
```

5.3.6 Analyse 6 : Calculer le nombre total de requêtes par code d'état :

```
grunt> total_requests_by_status = FOREACH (GROUP log_a BY Statucode) GENERATE group AS status_code, SUM(log_a.Statucode) AS total_requests;
grunt>
grunt>
grunt> dump total_requests_by_status;
```

```
(404,17318268)
(500,3185000)
(501,501)
(502,13055514)
(503,12451262)
(504,504)
```

5.3.7 Analyse des requêtes par pays

```
grunt> grouped_countries = GROUP log_a BY Country;
grunt> country_request_counts = FOREACH grouped_countries GENERATE group AS Country, COUNT(log_a) AS RequestCount;
grunt>
grunt>
grunt>
grunt>
```



```
(china,20803)
(india,13384)
(Canada,26940)
(Country,1)
(germany,19854)
(UNITED STATES,19019)
grunt>
```

6 Conclusion

L'analyse des logs d'un serveur web à l'aide d'Apache Pig offre une vision détaillée et précieuse du comportement des utilisateurs ainsi que des performances du site. En examinant les données telles que les adresses IP, les codes de statut HTTP, les navigateurs utilisés et les pays d'origine des visiteurs, nous pouvons tirer des enseignements précieux pour améliorer l'expérience utilisateur, optimiser les performances du site et prendre des décisions stratégiques.

À travers ce projet, nous avons pu mettre en œuvre différentes analyses telles que le comptage des requêtes par adresse IP, l'analyse des codes de statut HTTP, la répartition des navigateurs utilisés et l'analyse géographique des requêtes par pays. Ces analyses nous ont permis de mieux comprendre le comportement des utilisateurs, d'identifier les zones d'amélioration potentielles et de prendre des décisions éclairées pour optimiser le fonctionnement du site web.

En conclusion, l'utilisation d'Apache Pig pour l'analyse des logs de serveur web s'avère être un outil puissant pour extraire des informations pertinentes à partir de grandes quantités de données de manière efficace et éviter ainsi les erreurs et les incohérences. En continuant à explorer et à analyser ces données, les entreprises peuvent améliorer leur compréhension de leur audience en ligne et prendre des mesures pour optimiser leur présence sur le web.