

Διπλωματική Εργασία Αθανάσιος Ντόγαρης Α.Ε.Μ 53

Διαχείριση μεγάλου όγκου δεδομένων με το οικοσύστημα Hadoop και τη NoSQL βάση δεδομένων Hbase

“Big Data management with Hadoop ecosystem and HBase”

Ιούνιος 2018

Σκοπός της παρούσης εργασίας

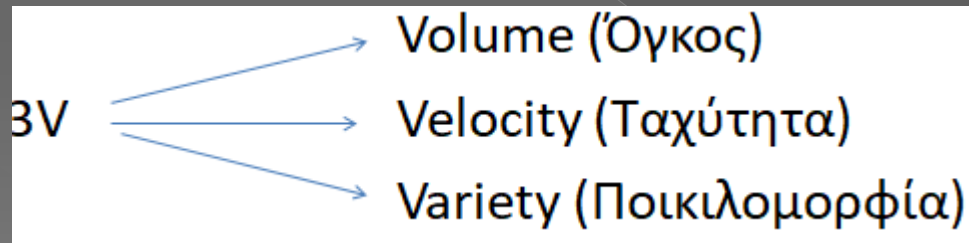
- Η μελέτη των δυνατοτήτων του οικοσυστήματος Hadoop και της HBase που είναι μία ανοιχτού λογισμικού NoSQL βάση δεδομένων ιδανική για αποθήκευση δεδομένων μεγάλου όγκου (Big Data).

Γιατί το Hadoop ήρθε στο προσκήνιο?

- Τα τελευταία χρόνια ο όγκος των δεδομένων έχει αυξηθεί εκθετικά διότι αυτά προέρχονται από διάφορες πηγές όπως τα Social Media, Κινητές συσκευές κ.α
- Αυτά τα δεδομένα τα οποία προέρχονται από πολλές και διαφορετικές πηγές ονομάζονται Δεδομένα Μεγάλου Όγκου (Big Data)

Ποιο ήταν το πρόβλημα?

- Οι παραδοσιακές βάσεις δεδομένων δεν είναι σχεδιασμένες για να χειρίζονται τα 3V των μεγάλων δεδομένων



Πως λύνει το Hadoop το πρόβλημα?

- Το Hadoop μπορεί να χειριστεί μεγάλου όγκου δεδομένα, μπορεί να επεξεργαστεί δεδομένα με μεγάλη ταχύτητα, και μπορεί να αποθηκεύσει οποιουδήποτε είδους δεδομένα.

Επομένως τι είναι το Hadoop?

- Το Hadoop ουσιαστικά είναι μία πλατφόρμα η οποία βασίζεται σε ένα κατακευμενένο σύστημα αρχείων το οποίο μας επιτρέπει να αποθηκεύσουμε δεδομένα πολύ μεγάλου όγκου σε διαφορετικούς υπολογιστές και μας παρέχει ένα API για την επεξεργασία αυτών των δεδομένων.
- Η βασική του ιδέα είναι ότι από τη στιγμή που τα δεδομένα αποθηκεύονται σε πολλούς ξεχωριστούς υπολογιστές μπορούμε να τα επεξεργαστούμε με έναν παράλληλο και κατακευμενένο τρόπο κατά τον οποίο ο κάθε υπολογιστής επεξεργάζεται τα δεδομένα τα οποία είναι αποθηκευμένα σε αυτόν.

Γιατί οι εταιρείες προτιμούν το Hadoop?

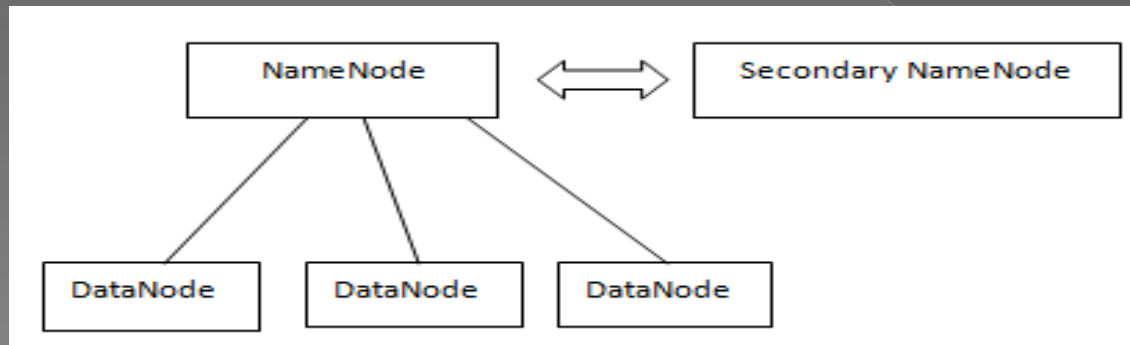
- 1) Γιατί είναι λογισμικό ανοιχτού κώδικα και δωρεάν
- 2) Γιατί μπορεί και χειρίζεται δεδομένα μεγάλου όγκου
- 3) Γιατί σχεδιάστηκε για να 'τρέχει' σε χαμηλού κόστους hardware

Ποια είναι τα κύρια συστατικά του Hadoop?

- Τα κύρια χαρακτηριστικά του Hadoop είναι:
- 1) Το HDFS
- 2) Το MapReduce
- 3) Το Yarn

Τι είναι το HDFS?

- Το HDFS (Hadoop Distributed File System) είναι ένα κατανεμημένο σύστημα αρχείων το οποίο μας επιτρέπει να αποθηκεύσουμε πολύ μεγάλα αρχεία σε διάφορους nodes ενός Hadoop cluster.



Τι είναι το MapReduce?

- Το MapReduce είναι μία προγραμματιστική πλατφόρμα η οποία μας επιτρέπει να υλοποιήσουμε παράλληλη και κατανεμημένη επεξεργασία σε δεδομένα μεγάλου όγκου σε ένα κατανεμημένο περιβάλλον όπως αυτό του Hadoop.
- Αποτελείται από δύο διεργασίες, την Map και την Reduce. Η Reduce διεργασία γίνεται μετά από την Map, διότι στην Reduce γίνεται η επεξεργασία των δεδομένων που προέρχονται από την Map.

Τι είναι το YARN?

- Το YARN (Yet Another Resource Negotiator) είναι ουσιαστικά η κεντρική μονάδα επεξεργασίας του Hadoop και είναι υπεύθυνο για την διαχείριση των διεργασιών και των πόρων του συστήματος

Τι είναι το Cloudera Quickstarts VM?

- Είναι μία πλατφόρμα χτισμένη σε Linux η οποία δημιουργήθηκε για επιχειρήσεις αλλά και για άτομα που θέλουν μία ολοκληρωμένη λύση για επεξεργασία μεγάλων δεδομένων.
- Περιέχει τα κύρια συστατικά του Hadoop (HDFS, MapReduce, YARN) καθώς και άλλα προγράμματα ανοιχτού λογισμικού όπως είναι η HBase και το Hive με τα οποία επίσης θα ασχοληθούμε.

Πλεονεκτήματα Cloudera VM

- Προσφέρει έναν πολύ απλό τρόπο για να ξεκινήσει κάποιος με το Hadoop. Όλα τα προγράμματα είναι προεγκατεστημένα.
- Έχει ένα απλό και διαχειρίσιμο interface.

Μειονεκτήματα Cloudera VM

- Έχει υψηλές απαιτήσεις συστήματος σε περίπτωση που εγκατασταθεί σε προσωπικό υπολογιστή:
 - 1) 64 – bit σύστημα
 - 2) Ελάχιστη μνήμη τουλάχιστον 4GB, αλλά στην πραγματικότητα χρειάζονται 8GB τουλάχιστο για επεξεργασία δεδομένων
 - 3) 4-πύρηνο επεξεργαστή
 - 4) Έρχεται με εγκαταστημένη την έκδοση 1.7 της Java η οποία χρειάζεται αναβάθμιση για να τρέξει το ElasticSearch

Προγράμματα που είναι εγκατεστημένα στο Cloudera VM και με τα οποία θα ασχοληθούμε

- ◉ Hue
- ◉ HBase
- ◉ Hive

Άλλα προγράμματα που είναι εγκατεστημένα στο Hadoop και δεν θα ασχοληθούμε είναι

- **Impala** – χρησιμοποιείται για τη δημιουργία ερωτημάτων SQL κατευθείαν στο HDFS και στην HBase
- **Pig** – προσφέρει μία γλώσσα υψηλού επιπέδου την PigLatin η οποία κάνει ποιο εύκολη την συγγραφή MapReduce κώδικα
- **Spark** – ένα ανοιχτού λογισμικού περιβάλλον για γρήγορη ανάλυση μεγάλων δεδομένων
- **Oozie** – ένα σύστημα προγραμματισμού ροών εργασιών
- **Solr** – Η αντίστοιχη μηχανή αναζήτησης που προσφέρεται με το Cloudera Quickstarts VM
- **Sqoop** – Επιτρέπει την αλληλεπίδραση των RDBMS με το HDFS

Τι είναι το Hue?

- Το Hue (Hadoop User Interface) παρέχει ένα web interface για το Hadoop και για τα προγράμματα που χρησιμοποιούνται για ανάλυση δεδομένων όπως είναι το Hive και η HBase.
- Επίσης με το Hue το ανέβασμα και η διαγραφή αρχείων από το HDFS γίνεται πιο εύκολα χωρίς τη χρήση της γραμμής εντολών

Τι είναι το Hive?

- Το Hive αρχικά είχε υλοποιηθεί από το Facebook και στη συνέχεια από την Apache Software Foundation. Μας παρέχει μία γλώσσα παρόμοια με την SQL, η οποία ονομάζεται Hive Query Language (ή αλλιώς HiveQL ή HQL) για τη δημιουργία ερωτημάτων πάνω στα δεδομένα που έχουν αποθηκευτεί σε μία Hadoop συστοιχία (cluster). Αυτό που κάνει είναι να μεταφράζει αυτά τα ερωτήματα σε MapReduce εργασίες, επιτρέποντας έτσι σε μη Java προγραμματιστές να επεξεργαστούν τα δεδομένα τους.
- Το Hive χρησιμοποιείται κυρίως σε εφαρμογές αποθηκών δεδομένων (data warehouse applications) όπου τα δεδομένα είναι σχεδόν στατικά και δεν αλλάζουν με γρήγορους ρυθμούς.

Τι είναι η HBase?

- Είναι η NoSQL πλατφόρμα του Hadoop γραμμένη σε Java
- Είναι column oriented δηλαδή οι τιμές μίας στήλης αποθηκεύονται όλες μαζί

Στην HBase οι πίνακες είναι τεσσάρων διαστάσεων. Σε αυτούς υπάρχουν:

- το κλειδί σειράς (row key)
- η οικογένεια στηλών (column family)
- η στήλη
- η χρονοσφραγίδα (timestamp)

Παράδειγμα πίνακα HBase

Εργασία				Προσωπικές Πληροφορίες		
<u>IdEργ</u>	Τμήμα	Βαθμός	Τίτλος	Όνομα	Διεύθυνση	ΑΜΚΑ

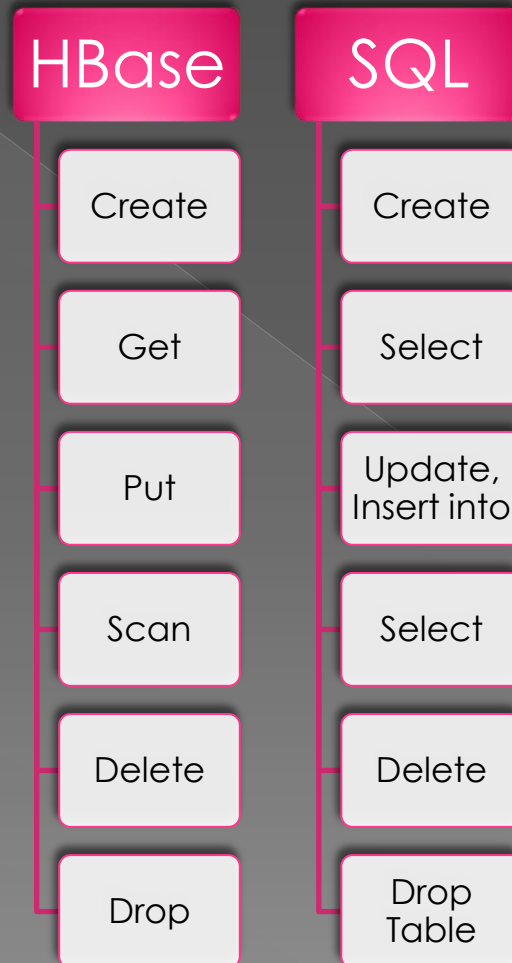
Σύγκριση πίνακα RDBMS με πίνακα της HBase

Id	Όνομα	Επίθετο	Τηλέφωνο	Πόλη
1	Γιώργος	Παπαδόπουλος	6999999999	Σέρρες
2	Γιάννης	Μακρυγιάννης	6999999998	Αθήνα
3	Φωτεινή	Ζησοπούλου	6999999997	Λάρισα
4	Νίκη	Δημητρίου	6999999996	Καβάλα
5	Αθανάσιος	Παυλίδης	6999999995	Σέρρες



Id	Στήλη	Τιμή
1	Όνομα	Γιώργος
1	Επίθετο	Παπαδόπουλος
1	Τηλέφωνο	6999999999
1	Πόλη	Σέρρες
2	Όνομα	Γιάννης
2	Επίθετο	Μακρυγιάννης
2	Τηλέφωνο	6999999998
2	Πόλη	Αθήνα
3	Όνομα	Φωτεινή
3	Επίθετο	Ζησοπούλου
3	Τηλέφωνο	6999999997
3	Πόλη	Λάρισα
4	Όνομα	Νίκη
4	Επίθετο	Δημητρίου
4	Τηλέφωνο	6999999996
4	Πόλη	Καβάλα
5	Όνομα	Αθανάσιος
5	Επίθετο	Παυλίδης
5	Τηλέφωνο	6999999995
5	Πόλη	Σέρρες

HBase εντολές vs SQL εντολές



HBase vs RDBMS

HBase

- Μέγεθος δεδομένων: Terrabytes, Petabytes
- Column oriented
- Ευέλικτο schema (μπορούμε να προσθέτουμε δεδομένα χωρίς περιορισμό)
- Εντολές API, εκτός κι αν συνδυαστεί με το Hive
- Ευρετήρια: Μόνο κλειδί σειράς, εκτός αν συνδυαστεί με το Hive ή το ElasticSearch
- Αρχιτεκτονική: Βασισμένη στο Hadoop, εύκολα διαβαθμίσιμη

SQL

- Μέγεθος δεδομένων: Gigabytes, Terrabytes.
- Row oriented (τις περισσότερες φορές)
- Σταθερό schema: ορίζεται κατά τον σχεδιασμό της βάσης
- SQL εντολές
- Ευρετήρια: Υπάρχουν
- Αρχιτεκτονική: Για πολύ μεγάλα δεδομένα χρειάζεται ακριβά συστήματα.

Τι είναι το Elasticsearch?

- Το Elasticsearch είναι μία ανοιχτού λογισμικού μηχανή αναζήτησης κειμένου καθώς και ανάλυσης δεδομένων η οποία μας επιτρέπει να αποθηκεύσουμε, να αναλύσουμε και να αναζητήσουμε πληροφορίες μέσα σε δεδομένα πολύ μεγάλου όγκου σε πολύ μικρό χρονικό διάστημα.
- Είναι γραμμένο στην γλώσσα Java και βασίζεται στο Apache Lucene. Είναι όμως πιο εύχρηστο, γιατί προσφέρει ένα απλό HTTP RESTful API το οποίο χρησιμοποιεί την JSON για την διαχείριση των δεδομένων. Επίσης με το Elasticsearch μπορεί να χρησιμοποιηθεί και το Kibana το οποίο είναι ένα εργαλείο ανάλυσης των δεδομένων με πολύ απλό τρόπο.

Παράδειγμα HDFS

File Browser

Search for file name ⚙️ Actions ▾ ✖ Move to trash ▾

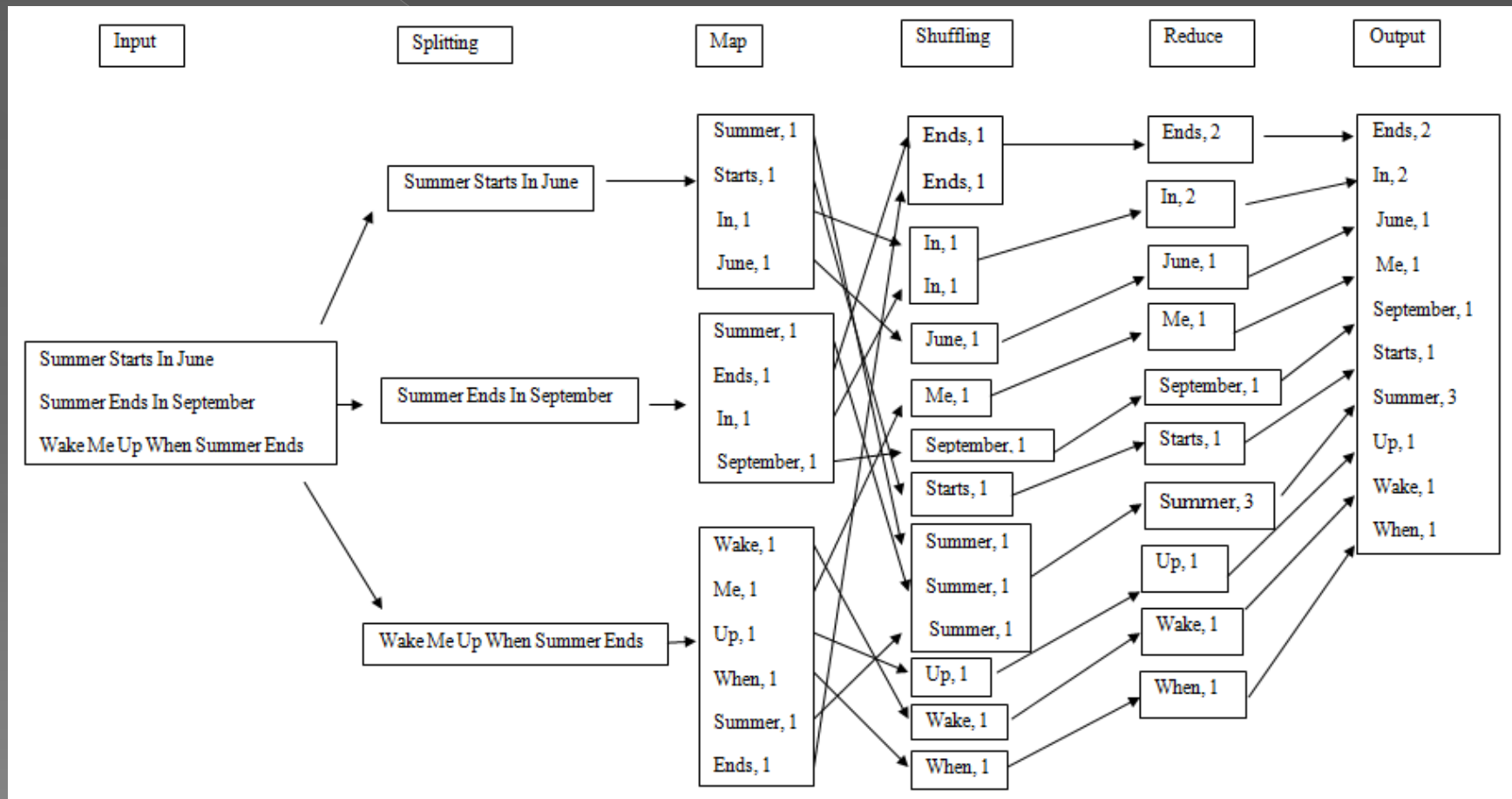
[Home](#) / [user](#) / [cloud](#)

<input type="checkbox"/>	Name	Size	User	Group
<input type="checkbox"/>	↑		cloudera	cloudera
<input type="checkbox"/>	.		cloudera	cloudera
<input checked="" type="checkbox"/>	Womens Clothing E-Commerce	3.1 MB	cloudera	cloudera

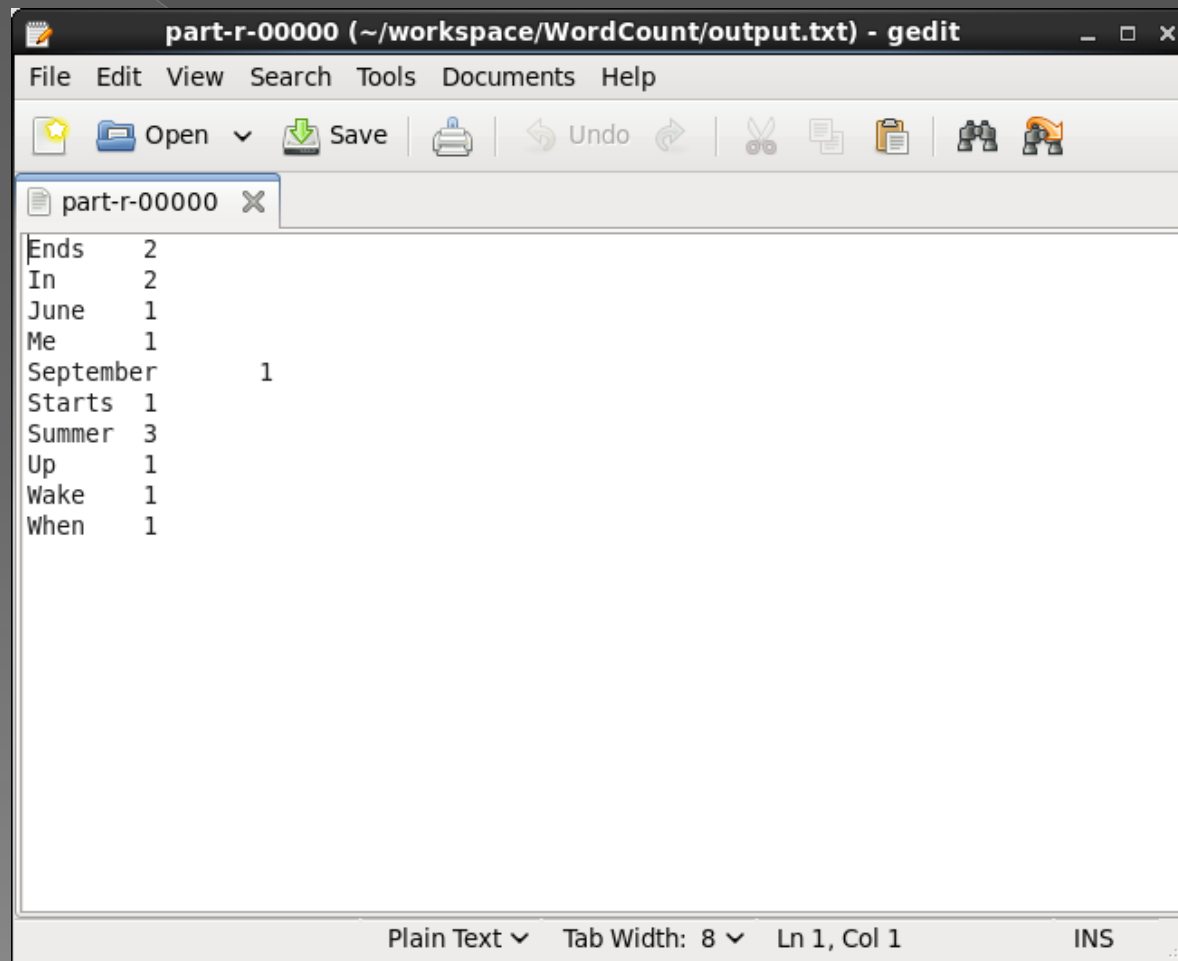
Show of 1 items Page

- Rename
- Move
- Copy
- Download
- Change permissions
- Summary
- Set replication

Παράδειγμα MapReduce



Παράδειγμα MapReduce



The screenshot shows a gedit window titled "part-r-00000 (~/.workspace/WordCount/output.txt) - gedit". The window contains a text file with the following content:

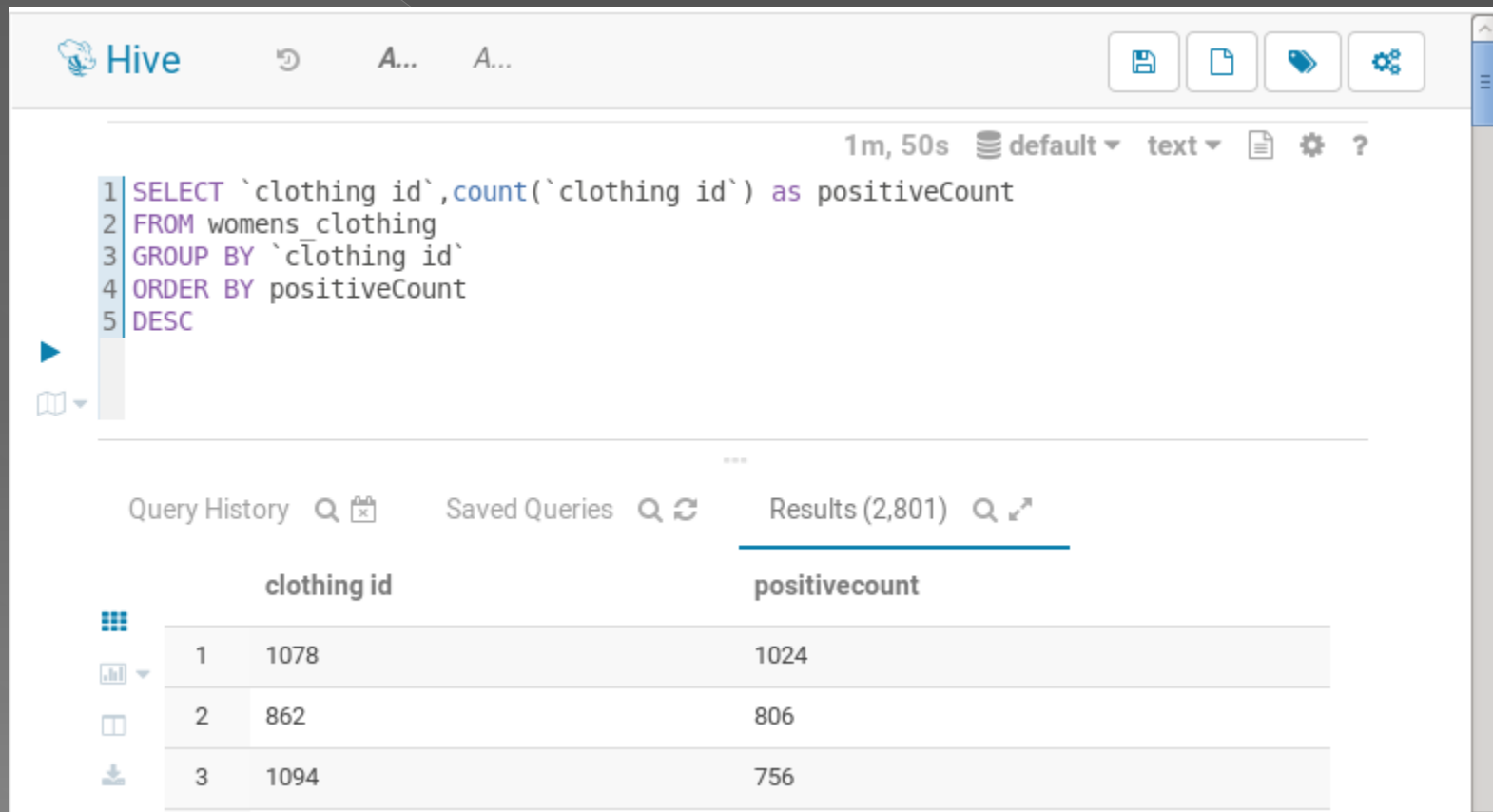
```
Ends 2
In 2
June 1
Me 1
September 1
Starts 1
Summer 3
Up 1
Wake 1
When 1
```

The status bar at the bottom of the window indicates "Plain Text", "Tab Width: 8", "Ln 1, Col 1", and "INS".

Παράδειγμα Hive

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
,Clothing	ID,Age,Title,Review Text,Rating,Recommended IND,Positive Feedback Count,Division Name,Department Name,Class Name																			
0,767,33,,	Absolutely wonderful - silky and sexy and comfortable,	4,1,0,Initmates,Intimate,Intimates																		
1,1080,34,,	"Love this dress! it's sooo pretty. i happened to find it in a store, and i'm glad i did bc i never would have ordered it online bc it's petite. i bought a petite and am 5'8". i love the length on me- hits just a li																			
2,1077,60,	Some major design flaws,"I had such high hopes for this dress and really wanted it to work for me. i initially ordered the petite small (my usual size) but i found this to be outrageously small. so small in fact																			
3,1049,50,	My favorite buy!,"I love, love, love this jumpsuit. it's fun, flirty, and fabulous! every time i wear it, i get nothing but great compliments!"	5,1,0,General Petite,Bottoms,Pants																		
4,847,47,	Flattering shirt,This shirt is very flattering to all due to the adjustable front tie. it is the perfect length to wear with leggings and it is sleeveless so it pairs well with any cardigan. love this shirt!!!	5,1,6,General																		
5,1080,49,	Not for the very petite,"I love tracy reese dresses, but this one is not for the very petite. i am just under 5 feet tall and usually wear a 0p in this brand. this dress was very pretty out of the package but its a l																			
6,858,39,	Cagrc coal shimmer fun,"I aded this in my basket at hte last mintue to see what it would look like in person. (store pick up). i went with teh darkler color only because i am so pale :-) hte color is really gorgeou																			
7,858,39,	"Shimmer, surprisingly goes with lots", "I ordered this in carbon for store pick up, and had a ton of stuff (as always) to try on and used this top to pair (skirts and pants). everything went with it. the color is rea																			
8,1077,24,	Flattering,I love this dress. i usually get an xs but it runs a little snug in bust so i ordered up a size. very flattering and feminine with the usual retailer flair for style.,	5,1,0,General,Dresses,Dresses																		
9,1077,34,	Such a fun dress!,"I'm 5'5" and 125 lbs. i ordered the s petite to make sure the length wasn't too long. i typically wear an xs regular in retailer dresses. if you're less busty (34b cup or smaller), a s petite will																			
10,1077,53,	Dress looks like it's made of cheap material,Dress runs small esp where the zipper area runs. i ordered the sp which typically fits me and it was very tight! the material on the top looks and feels very cheap																			
11,1095,39,	This dress is perfection! so pretty and flattering.,	5,1,2,General Petite,Dresses,Dresses																		

Παράδειγμα Hive



The screenshot shows the Hive web interface. At the top, there is a navigation bar with the Hive logo, a refresh icon, and two 'A...' buttons. On the right, there are icons for save, print, share, and settings. Below the navigation bar, the execution time is shown as '1m, 50s' along with a 'default' dropdown and a 'text' dropdown. The main area contains a SQL query:

```
1 SELECT `clothing id`,count(`clothing id`) as positiveCount
2 FROM womens_clothing
3 GROUP BY `clothing id`
4 ORDER BY positiveCount
5 DESC
```

Below the query, there are tabs for 'Query History', 'Saved Queries', and 'Results (2,801)'. The 'Results' tab is active, showing a table with two columns: 'clothing id' and 'positivecount'. The table contains three rows of data:

	clothing id	positivecount
1	1078	1024
2	862	806
3	1094	756

Παράδειγμα: Εντολές HBase

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
{NAME => 'professional data', DATA_BLOCK_ENCODING => 'NONE', BLOOMFILTER => 'ROW',  
, REPLICATION_SCOPE => '0', VERSIONS => '1', COMPRESSION => 'NONE', MIN_VERSION  
S => '0', TTL => 'FOREVER', KEEP_DELETED_CELLS => 'FALSE', BLOCKSIZE => '65536',  
IN_MEMORY => 'false', BLOCKCACHE => 'true'}  
2 row(s) in 0.1450 seconds  
  
hbase(main):005:0> count 'employee'  
0 row(s) in 0.0750 seconds  
  
=> 0  
hbase(main):006:0> put 'employee',1,'personal data:name','George Pappas'  
0 row(s) in 0.1120 seconds  
put 'employee',1,'personal data:marital_status','unmarried'  
0 row(s) in 0.0080 seconds  
  
hbase(main):008:0> scan 'employee'  
ROW COLUMN+CELL  
1 column=personal data:marital_status, timestamp=15200097715  
70, value=unmarried  
1 column=personal data:name, timestamp=1520008119616, value=  
George Pappas  
1 row(s) in 0.0250 seconds  
  
hbase(main):009:0> █
```

Παράδειγμα HBase μέσω Hue

Home - HBase / Employee Switch Cluster ▾

row_key, row_prefix* +scan_len [col1, family:col2, fam3:, col_prefix* +3, professional_data: personal_data:

All Sort By ASC ▾

1	professional_data: Salary	personal_data: Lastname	personal_data: Firstname	professional_data: Department	personal_data: status	professional_data: Jobtitle
	2300	Papagiannis	George	IT	married	Web Developer
2	professional_data: Salary	personal_data: Lastname	personal_data: Firstname	professional_data: Department	personal_data: status	professional_data: Jobtitle
	1500	Makris	Dimitris	Logistics	Divorced	Data Entry

9 seconds.

3

Δημιουργία ευρετηρίου ElasticSearch - Kibana

Management / Kibana

Index Patterns Saved Objects Advanced Settings

+ Create Index Pattern

★ student

★ student ★ ↻ 🗑️

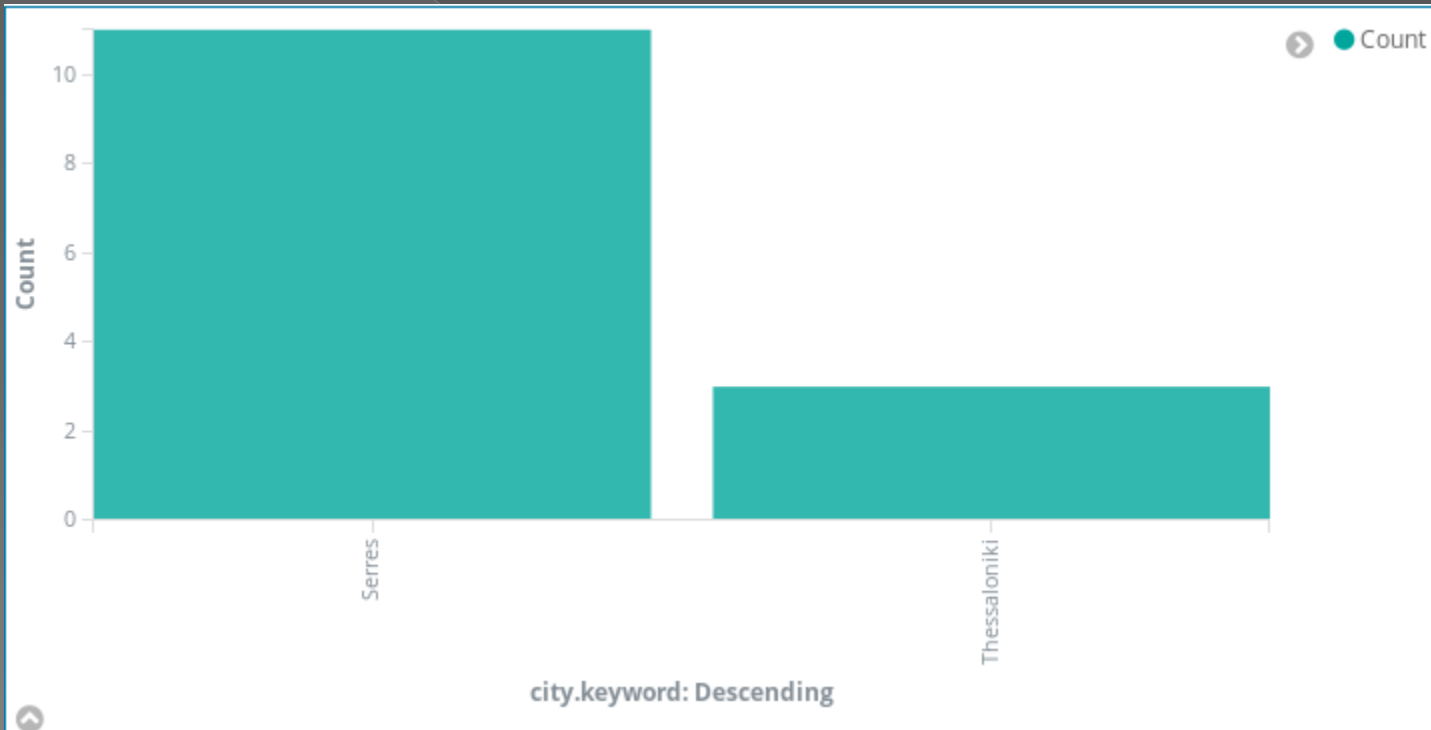
This page lists every field in the **student** index and the field's associated core type as recorded by Elasticsearch. While this list allows you to view the core type of each field, changing field types must be done using Elasticsearch's [Mapping API](#).

fields (14) scripted fields (0) source filters (0)

Filter All field types

name	type	format	searchable	aggregatable	excluded	controls
_id	string		✓	✓		
_index	string		✓	✓		
_score	number					
_source	_source					
_type	string		✓	✓		
city	string		✓			
city.keyword	string		✓	✓		
country	string		✓			
country.keyword	string		✓	✓		
firstname	string		✓			
firstname.keyword	string		✓	✓		
lastname	string		✓			
lastname.keyword	string		✓	✓		
stid	number		✓	✓		

Οπτικοποίηση δεδομένων με Kibana



Συμπεράσματα

- Παρουσιάστηκε το **HDFS** και η ευκολία να ανεβάσουμε κάποιο αρχείο σε αυτό μέσω του **Hue**
- Παρουσιάστηκε το **MapReduce** και πως μπορούμε να τρέξουμε ένα απλό **WordCount** πρόγραμμα
- Παρουσιάστηκε το **Hive** και πως μπορούμε να θέσουμε ερωτήματα σε μία βάση μέσω του **Hue**
- Παρουσιάστηκαν οι εντολές της **HBase Shell** καθώς και η δημιουργία πινάκων στην **HBase** μέσω του **Hue**
- Τέλος δημιουργήθηκαν ευρετήρια στα δεδομένα μας με τη βοήθεια του **ElasticSearch** και κατ'επέκταση του **Kibana**

Συμπεράσματα

Ο συνδυασμός του Elasticsearch με την HBase (παρόλο που το Elasticsearch μπορεί να χρησιμοποιηθεί κι αυτό ως βάση δεδομένων) μπορεί να προσφέρει περισσότερη ασφάλεια, ώστε να μη χαθεί κάποιο μέρος ενός αρχείου, καθώς και μεγαλύτερη ταχύτητα στην αναζήτηση και οπτικοποίηση των δεδομένων μας.

Αναφορές

- ◉ Ιστοσελίδες
- ◉ [1] WordCount Example on CDH 5.12 διαθέσιμο online: <https://www.youtube.com/watch?v=3l7gAHjVOc4> [πρόσβαση 17/06/2018]
- ◉ [2] MapReduce Tutorial διαθέσιμο online: https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html#Example: WordCount_v1.0 [πρόσβαση 17/06/2018]
- ◉ [3] What is Big Data <https://www.edureka.co/blog/hadoop-tutorial/#WhatsBigData> [πρόσβαση 17/06/2018]
- ◉ [4] Big Data Tutorial <https://www.edureka.co/blog/big-data-tutorial> [πρόσβαση 17/06/2018]
- ◉ [6] HDFS Tutorial <https://www.edureka.co/blog/hdfs-tutorial> [πρόσβαση 17/06/2018]
- ◉ [7] HDFS Architecture <https://www.edureka.co/blog/apache-hadoop-hdfs-architecture/> [πρόσβαση 17/06/2018]
- ◉ [8] MapReduce Tutorial <https://www.edureka.co/blog/mapreduce-tutorial/> [πρόσβαση 17/06/2018]
- ◉ [9] Getting started with HBase: The Hadoop database <https://www.pluralsight.com/courses/hbase-hadoop-getting-started> [πρόσβαση 17/06/2018]
- ◉ [10] HBase Architecture: HBase Data Model & HBase Read/Write Mechanism <https://www.edureka.co/blog/hbase-architecture/> [πρόσβαση 17/06/2018]
- ◉ [11] ElasticSearch Reference: Getting Started https://www.elastic.co/guide/en/elasticsearch/reference/current/basic_concepts.html [πρόσβαση 17/06/2018]
- ◉ Βιβλία
- ◉ [5] Aven Jeffrey (2017), Sams Teach Yourself Hadoop in 24 Hours, Εκδόσεις Sams Publishing