

Βάσεις Δεδομένων II
Εργαστηριακή Άσκηση 2020/21

Όνομα	Επώνυμο	ΑΜ
Νικηφόρος – Γεώργιος	Παπαγεωργίου	1059633
Νικόλαος	Σταμόπουλος	1057764

Βεβαιώνω ότι είμαι συγγραφέας της παρούσας εργασίας και ότι έχω αναφέρει ή παραπέμψει σε αυτήν, ρητά και συγκεκριμένα, όλες τις πηγές από τις οποίες έκανα χρήση δεδομένων, ιδεών, προτάσεων ή λέξεων, είτε αυτές μεταφέρονται επακριβώς (στο πρωτότυπο ή μεταφρασμένες) είτε παραφρασμένες. Επίσης βεβαιώνω ότι αυτή η εργασία προετοιμάστηκε από εμένα προσωπικά ειδικά για το συγκεκριμένο μάθημα/σεμινάριο/πρόγραμμα σπουδών.

Έχω ενημερωθεί ότι σύμφωνα με τον εσωτερικό κανονισμό λειτουργίας του Πανεπιστημίου Πατρών άρθρο 50§6, τυχόν προσπάθεια αντιγραφής ή εν γένει φαλκίδευσης της εξεταστικής και εκπαιδευτικής διαδικασίας από οιονδήποτε εξεταζόμενο, πέραν του μηδενισμού, συνιστά βαρύ πειθαρχικό παράπτωμα.

Υπογραφή



22 / 09 / 2021

Υπογραφή



22 / 09 / 2021

Συνημμένα αρχεία κώδικα

Μαζί με την παρούσα αναφορά υποβάλλουμε τα παρακάτω αρχεία κώδικα

Αρχείο	Αφορά το ερώτημα	Περιγραφή/Σχόλιο
queryX.ipynb	1	Περιέχει τον κώδικα για το ερώτημα X, όπου $X \in [1, 10]$

Τεχνικά χαρακτηριστικά περιβάλλοντος λειτουργίας

Τεχνικά χαρακτηριστικά φυσικού Η/Υ που χρησιμοποιήθηκε για την εργασία

Χαρακτηριστικό	Τιμή
CPU model	AMD Ryzen 5 2600 Six-Core Processor
CPU clock speed	3.4GHz
Physical CPU cores	6
Logical CPU cores	12
RAM	16Gb
Secondary Storage Type	SSD

Τεχνικά χαρακτηριστικά εικονικής μηχανής (VM) που χρησιμοποιήθηκε για την εργασία

Χαρακτηριστικό	Τιμή
CPU cores	6
Execution cap	100%
RAM	9Gb
VM OS	Ubuntu 20.04.2 LTS
VM software	VirtualBox
Host OS	Windows 10

Ερώτημα 1: Απαντήσεις ερωτημάτων

Σημείωση: Τα ζητούμενα αποτελέσματα παρατίθενται σε μορφή screenshot, για την καλύτερη απεικόνισή τους.

Ερώτημα	Απάντηση
Δώστε το πλήθος των χρηστών που είδαν την ταινία "Jumanji".	<pre>+-----+-----+ title total_viewers +-----+-----+ Jumanji (1995) 22243 +-----+-----+</pre>
Δώστε τα ονόματα των ταινιών που οι χρήστες χαρακτήρισαν ως "boring".	<pre>+-----+-----+-----+ title lower_tag +-----+-----+-----+ (500) Days of Summer (2009) boring 101 Reykjavik (101 Reykjavík) (2000) boring 12 Years a Slave (2013) boring 1408 (2007) boring 1492: Conquest of Paradise (1992) boring +-----+-----+-----+</pre>
Δώστε τους χρήστες που έχουν χαρακτηρίσει την ταινία ως "Bollywood" και την έχουν αξιολογήσει με βαθμό >3.	<pre>+-----+-----+-----+ userId rating lower_tag +-----+-----+-----+ 10573 4.0 bollywood 19837 5.0 bollywood 23333 4.0 bollywood 25004 5.0 bollywood 31338 4.5 bollywood +-----+-----+-----+</pre>

Βρείτε τις 10 κορυφαίες ταινίες για κάθε έτος.	<div> <div>Before the Fall (NaPolA - Elite für den Führer) (2004)</div> <div>2005</div> <div>5.0</div> <div>1</div> </div> <div> <div>Dancemaker (1998)</div> <div>2005</div> <div>5.0</div> <div>2</div> </div> <div> <div>Fear Strikes Out (1957)</div> <div>2005</div> <div>5.0</div> <div>3</div> </div> <div> <div>Gate of Heavenly Peace, The (1995)</div> <div>2005</div> <div>5.0</div> <div>4</div> </div> <div> <div>Life Is Rosy (a.k.a. Life Is Beautiful) (Vie est belle, La) (1987)</div> <div>2005</div> <div>5.0</div> <div>5</div> </div> <div> <div>Married to It (1991)</div> <div>2005</div> <div>5.0</div> <div>6</div> </div> <div> <div>My Life and Times With Antonin Artaud (En compagnie d'Antonin Artaud) (1993)</div> <div>2005</div> <div>5.0</div> <div>7</div> </div> <div> <div>Not Love, Just Frenzy (Más que amor, frenesí) (1996)</div> <div>2005</div> <div>5.0</div> <div>8</div> </div> <div> <div>Paris Was a Woman (1995)</div> <div>2005</div> <div>5.0</div> <div>9</div> </div> <div> <div>Take Care of My Cat (Goyangileul butaghae) (2001)</div> <div>2005</div> <div>5.0</div> <div>10</div> </div>
Δώστε τις ετικέτες για κάθε ταινία και το όνομα της ταινίας για το έτος 2015.	<div> <div>"Great Performances" Cats (1998)</div> <div>[[BD-R]</div> </div> <div> <div>'burbs, The (1989)</div> <div>[[1980's, black comedy, dark comedy, Joe Dante, quirky]</div> </div> <div> <div>(500) Days of Summer (2009)</div> <div>[[annoying, artistic, bad dialogue, boring, depressing, Joseph Gordon-Levitt, overrated, slow, stupid, Zooey Deschanel, intelligent, nonlinear, artistic, bittersweet, Funny, humor, humorous, intelligent, Joseph Gordon-Levitt, music, nonlinear, quirky, relationships, romance, Zooey Deschanel, bittersweet, quirky, romance, Joseph Gordon-Levitt, artistic, no happy ending, nonlinear, overrated]]</div> </div> <div> <div>...tick... tick... tick... (1970)</div> <div>[[BD-R]</div> </div> <div> <div>l (2014)</div> <div>[[Sukumar]</div> </div>
Δώστε το πλήθος των ratings για κάθε ταινία.	<div> <div>title</div> <div>total_ratings</div> </div> <div> <div>Pulp Fiction (1994)</div> <div>67310</div> </div> <div> <div>Forrest Gump (1994)</div> <div>66172</div> </div> <div> <div>Shawshank Redemption, The (1994)</div> <div>63366</div> </div> <div> <div>Silence of the Lambs, The (1991)</div> <div>63299</div> </div> <div> <div>Jurassic Park (1993)</div> <div>59715</div> </div>
Βρείτε τους 10 πρώτους χρήστες με τα περισσότερα rating για κάθε χρονιά.	<div> <div>userId</div> <div>yearNum</div> <div>total_ratings</div> <div>rank</div> </div> <div> <div>131160</div> <div>1995</div> <div>3</div> <div>1</div> </div> <div> <div>28507</div> <div>1995</div> <div>1</div> <div>2</div> </div>
Βρείτε τις ταινίες με τα περισσότερα ratings για κάθε κατηγορία ταινίας.	<div> <div>genres</div> <div>title</div> <div>total_ratings</div> </div> <div> <div>(no genres listed)</div> <div>Doctor Who: The Time of the Doctor (2013)</div> <div>36</div> </div> <div> <div>Action</div> <div>Jurassic Park (1993)</div> <div>59715</div> </div> <div> <div>Adventure</div> <div>Jurassic Park (1993)</div> <div>59715</div> </div> <div> <div>Animation</div> <div>Toy Story (1995)</div> <div>49695</div> </div> <div> <div>Children</div> <div>Toy Story (1995)</div> <div>49695</div> </div>
Δώστε το σύνολο των χρηστών που παρακολουθούν την ίδια ταινία, την ίδια μέρα και ώρα.	<div> <div>total_viewers</div> </div> <div> <div>4281178</div> </div>
Δώστε το πλήθος των ταινιών, για κάθε κατηγορία, που οι χρήστες χαρακτήρισαν ως "funny" και με rating > 3.5.	<div> <div>genres</div> <div>movies_count</div> </div> <div> <div>Action</div> <div>431</div> </div> <div> <div>Adventure</div> <div>465</div> </div> <div> <div>Animation</div> <div>268</div> </div> <div> <div>Children</div> <div>273</div> </div> <div> <div>Comedy</div> <div>1618</div> </div>

Ερώτημα 2: Σύγκριση επιδόσεων σε single node/virtual cluster/Livy

Ρυθμίσεις virtual cluster

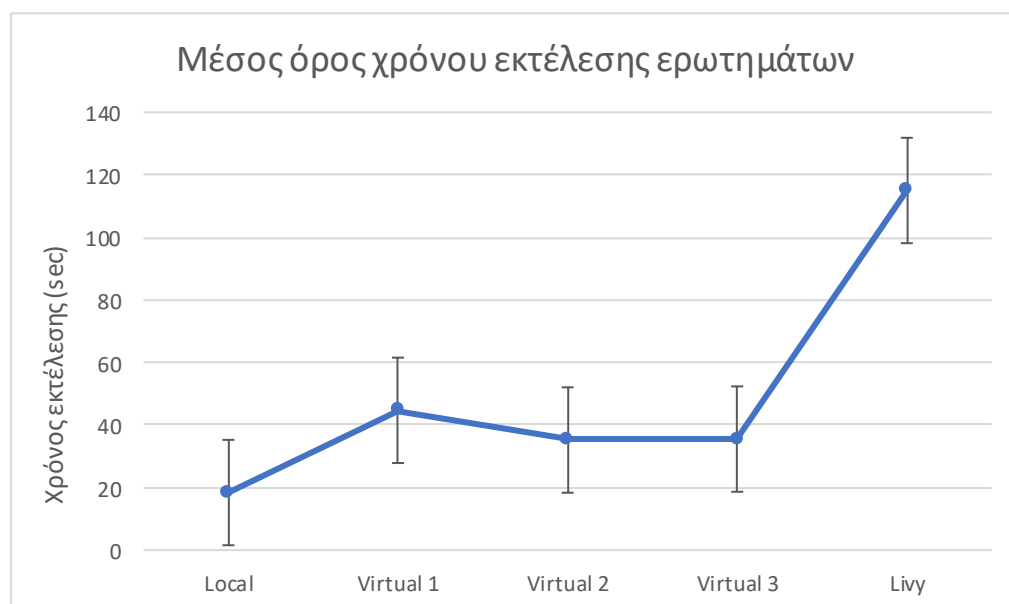
A/A	Executor cores	Executor mem	Driver cores	Driver mem
1	1	1G	1	1G
2	2	2G	1	1G
3	2	2G	2	2G

Χρόνοι εκτέλεσης

Σημείωση: Οι παρακάτω χρόνοι εκτέλεσης μετρήθηκαν σε δευτερόλεπτα (*seconds*). Η χρονομέτρηση έγινε με χρήση της βιβλιοθήκης *sparkMeasure*.

Ερώτημα	Local	Virtual 1	Virtual 2	Virtual 3	Livy
1	14	40	28	29	120
2	5	19	15	15	15
3	20	46	40	38	123
4	26	55	44	52	180
5	4	18	15	13	16
6	15	40	36	34	108
7	27	53	41	41	168
8	25	48	41	38	119
9	32	90	62	63	181
10	16	38	30	32	120

Ανάλυση αποτελεσμάτων



Μετά την χρονομέτρηση και ανάλυση των αποτελεσμάτων εξάγουμε τις παρακάτω παρατηρήσεις:

- Σε single node μηχανήμα (local) πετυχαίνουμε τάχιστα εκτέλεση των ερωτημάτων, καθώς χρησιμοποιείται όλη η υπολογιστική δύναμη που έχουμε αναθέσει στο VM.
- Ανάμεσα στα Virtual 1, Virtual 2 και Virtual 3, παρατηρούμε πως το πρώτο καταναλώνει περισσότερο χρόνο για την εκτέλεση των ερωτημάτων σε σχέση με τα άλλα δύο και αυτό λογικά οφείλεται στην ανάθεση μόνο ενός πυρήνα (core) για κάθε worker. Επίσης παρατηρούμε πως ανάμεσα στα Virtual 2 και 3 οι διαφορές είναι μηδαμινές, συνεπώς η αύξηση των πυρήνων και της μνήμης που δεσμεύει ο driver στο Virtual 3 δεν έφερε καλύτερα αποτελέσματα από το 2.
- Η εκτέλεση των ερωτημάτων στον Livy server καταναλώνει τον μέγιστο χρόνο.
- Όσον αφορά τα ερωτήματα, παρατηρούμε πως συγκεκριμένα τα ερωτήματα 2 και 5 τρέχουν πάντα σε πολύ λιγότερο χρόνο σε σύγκριση με τα υπόλοιπα. Πιθανολογούμε πως αυτό οφείλεται στο γεγονός ότι αυτά τα δύο ερωτήματα δεν κάνουν χρήση του αρχείου rating.csv, το οποίο περιέχει το μεγαλύτερο πλήθος εγγραφών απ' όλα και συνεπώς είναι πιο «ακριβό» σε πλήθος πράξεων.

Βιβλιογραφία

1. *PySpark 3.1.2 Documentation*. <http://spark.apache.org/docs/latest/api/python/>
2. Α. Κομνηνός. *Φροντιστήριο 6 – Εισαγωγή στο Apache Spark*.
<https://eclass.upatras.gr/modules/document/index.php?course=CEID1176&openDir=/5e6f65ear83d/60756472Ab51>
3. Nishant Bahri. *Movie Lens Data Analysis Using PySpark [for beginners]*.
<https://medium.com/analytics-vidhya/movie-lens-data-analysis-using-pyspark-for-beginners-9c0f5f21eaf5>
4. Mauro Krikorian. *Movie Data Statistics with Apache Spark*.
<https://medium.com/southworks/movie-data-statistics-with-apache-spark-58c2ef8fe452>