TDT4305 Big Data Project

Phase 1 Data Analysis with Spark

Due Date: 14 March 23:59

Project Assistance:

Stella Maropaki - stella.maropaki@ntnu.no Shweta Tiwari - shweta.tiwari@ntnu.no Özer Özdikis - ozer.ozdikis@ntnu.no

General information

- Consultation hours in room 248 IT-bygget:
 - Wednesday 10-12
 - Friday 10-12



Deadline phase 2

Objectives









Music Data

• albums.csv

id	artist_id	album_title	genre	year_of_pub	num_of_tracks	num_of_sales	rolling_stone_critic	mtv_critic	music_maniac_critic
1	1767	Call me Cat Moneyless That Doggies	Folk	2006	11	905193	4	1.5	3
2	23548	Down Mare	Metal	2014	7	969122	3	4	5
3	17822	Embarrassed Hungry	Latino	2000	11	522095	2.5	1	2
4	19565	Standard Immediate Engineer Slovakia	Pon	2017	4	610116	15	2	4

artists.csv

id	real_name	art_name	role	year_of_birth	country	city	email	zip_code
1	Kameko Nelson		female voice	1981	Slovenia	Vedrin	lorem.tristique.aliquet@nonummyFusce.co.uk	6914
2	Sacha Cameron	Bianco Ibureno Chloride	artist	1955	Guernsey	Montigny-le-Tilleul	Sed@elementum.com	0568
3	Thane E. Oliver		rapper	1993	Guinea-Bissau	Saint-Sébastien-sur-Loire	Integer.tincidunt.aliquam@libero.com	813154
4	Cole E. Joseph	Persian responsible	bassist	1994	Estonia	Newauav	nulla.magna.malesuada@vulputate.com	A4S 2B3

Notes

- Use Spark functions:
 - map(), reduceByKey(), sortByKey(), reduce(), count(), ...
 - no external libraries
- For efficiency:
 - use only data you need
- For faster testing:
 - o rdd_sample = rdd.sample(false, 0.1, 5)
- Plagiarism:
 - Use private repositories (Git, BitBucket, ...)
 - Don't copy

Notes

Part of the evaluation of the project will be done by running the code files with automatic script.

Use the requested:

- Code file names
- Output formats
- Output file names

Make sure your code is running in terminal.

Deliver on Blackboard

- Code files:
 - 9 tasks using RDD API
 - 1 task using Dataset API
- Report (pdf):
 - short description of code for each task
 - results for needed tasks
- Output files (.tsv):
 - for the tasks needed
 - in the described format
- Do not upload the datasets!!

The following expected results are not the correct ones!

They only have the expected format!

How many genres are there?

id	artist_id	album_title	genre	year_of_pub	num_of_tracks	num_of_sales	rolling_stone_critic	mtv_critic	music_maniac_critic
1	1767	Call me Cat Moneyless That Doggies	Folk	2006	11	905193	4	1.5	3
2	23548	Down Mare	Metal	2014	7	969122	3	4	5
3	17822	Embarrassed Hungry	Latino	2000	11	522095	2.5	1	2
4	19565	Standard Immediate Engineer Slovakia	Pon	2017	4	610116	1.5	2	4

Expected format of result: 25

What is the year of birth of the oldest artist?

id	real name	art name	role	year of birth	country	city	email	zip_code
1	Kameko Nelson		female voice	1981	Slovenia	Vedrin	lorem.tristique.aliquet@nonummyFusce.co.uk	6914
2	Sacha Cameron	Bianco Ibureno Chloride	artist	1955	Guernsey	Montigny-le-Tilleul	Sed@elementum.com	0568
3	Thane E. Oliver		rapper	1993	Guinea-Bissau	Saint-Sébastien-sur-Loire	Integer.tincidunt.aliquam@libero.com	813154
4	Cole E. Joseph	Persian responsible	bassist	1994	stonia	Newauav	nulla.magna.malesuada@vulputate.com	A4S 2B3

Expected format of results: 1955

Find the total number of artists coming from each country and sort them in descending order of artists counts.

id	real_name	art_name	role	year_of_birth	country	tity	email	zip_code
1	Kameko Nelson		female voice	1981	Slovenia	Vedrin (lorem.tristique.aliquet@nonummyFusce.co.uk	6914
2	Sacha Cameron	Bianco Ibureno Chloride	artist	1955	Guernsey	Nontigny-le-Tilleul	Sed@elementum.com	0568
3	Thane E. Oliver		rapper	1993	Guinea-Bissau	Saint-Sébastien-sur-Loire	Integer.tincidunt.aliquam@libero.com	813154
4	Cole E. Joseph	Persian responsible	bassist	1994	Estonia	lewauav	nulla.magna.malesuada@vulputate.com	A4S 2B3

Expected format of results:

Botswana 156

China 156

Armenia 130

Find the total number of albums each artist has and sort them in descending order of album counts.

id	artist_id	album_title	genre	year_of_pub	num_of_tracks	num_of_sales	rolling_stone_critic	mtv_critic	music_maniac_critic
1	1767	Call me Cat Moneyless That Doggies	Folk	2006	11	905193	4	1.5	3
2	23548	Down Mare	Metal	2014	7	969122	3	4	5
3	17822	Embarrassed Hungry	Latino	2000	11	522095	2.5	1	2
4	19565	Standard Immediate Engineer Slovakia	Pon	2017	4	610116	15	2	4

Expected format of results:

10708 11

1608 9

47714 9

Find the total number of sales per genre and sort them in descending order of album sales.

id	artist_id	album_title	genre	year_of_pub	num_of_tracks	num_of_sales	rolling_stone_critic	mtv_critic	music_maniac_critic
1	1767	Call me Cat Moneyless That Doggies	Folk	2006	11	905193	4	1.5	3
2	23548	Down Mare	Metal	2014	7	969122	3	4	5
3	17822	Embarrassed Hungry	Latino	2000	11	522095	2.5	1	2
4	19565	Standard Immediate Engineer Slovakia	Pon	2017	4	610116	1.5	2	4

Expected format of results:

Indie 23848059

Pop 17166379

Rap 8889141

Find the top 10 albums with the best average critic.

id	artist id	album title	genre	year of pub	num of tracks	num of sales	rolling stone critic	mtv critic	music maniac critic
1	1767	Call me Cat Moneyless That Doggies	Folk	2006	11	905193	4	1.5	3
2	23548	Down Mare	Metal	2014	7	969122	3	4	5
3	17822	Embarrassed Hungry	Latino	2000	11	522095	2.5	1	2
4	19565	Standard Immediate Engineer Slovakia	Pon	2017	Δ	610116	15	2	4

Expected format of results:

156 5.0

189 5.0

267 4.9

For the 10 albums of task 6, find the countries of their artists.

id	artist_id	album_title	genre	year_of_pub	num_of_tracks	num_of_sales	rolling_stone_critic	mtv_critic	nusic_maniac_critic
1	1767	Call me Cat Moneyless That Doggies	Folk	2006	11	905193	4	1.5	1
2	23548	Down Mare	Metal	2014	7	969122	3	4	5
3	17822	Embarrassed Hungry	Latino	2000	11	522095	2.5	1	2
4	19565/	Standard Immediate Engineer Slovakia	Pon	2017	4	610116	1.5	2	

id	real_name	art_name	role	year_of_birth	country	city	email	zip_code
1	Kameko Nelson		female voice	1981	Slovenia	Vedrin	lorem.tristique.aliquet@nonummyFusce.co.uk	6914
2	Sacha Cameron	Bianco Ibureno Chloride	artist	1955	Guernsey	Montigny-le-Tilleul	Sed@elementum.com	0568
3	Thane E. Oliver		rapper	1993	Guinea-Bissau	Saint-Sébastien-sur-Loire	Integer.tincidunt.aliquam@libero.com	813154
4	Cole E. Joseph	Persian responsible	bassist	1994	Estonia	Newquav	nulla.magna.malesuada@vulputate.com	A4S 2B3

Expected format of results:

156 5.0 Slovenia

189 5.0 Guernsey

267 4.9 Estonia

Find the artists that have an album with the highest (5.0) MTV critic, sorted alphabetically.

id	artist_id	album_titl	e	genre	year_of	_pub num_	of_tracks	num_of_sales	rolling_stone_critic	mtv_critic r	nusic_maniac_critic
1	1767	Call me Ca	at Moneyless That Doggies	Folk	2006	11		905193	4	1.5	3
2	23548	Down Mare	9	Metal	2014	7		969122	3	4 5	5
3	17822	Embarrass	ed Hungry	Latino	2000	11		522095	2.5	1 2	2
4	19565	Standard I	mmediate Endineer Slovakia	Pon	2017	4		610116	15	2	ſ
id	real_nam	ne ne	art_name	role	year_of_birth	country	city		email		zip_code
1	Kameko I	Nelson		female voice	1981	Slovenia	Vedrin		lorem.tristique.alique	t@nonummyFusce.d	co.uk 6914
2	Sacha Ca	ameron	Bianco Ibureno Chloride	artist	1955	Guernsey	Montigr	ny-le-Tilleul	Sed@elementum.cor	m	0568

Guinea-Bissau

Estonia

Saint-Sébastien-sur-Loire

Newquay

Integer.tincidunt.aliquam@libero.com

nulla.magna.malesuada@vulputate.com

813154

A4S 2B3

1993

1994

rapper

bassist

Expected format of results:

Persian responsible

Sacha Cameron

Thane E. Oliver

Cole E. Joseph

Thane E. Oliver

Cole E. Joseph

Find the average MTV critic of all albums for each artist from Norway (country='Norway').

											1000			
id	artist_id	album_title			genre	у	ear_of_p	oub num_of_	tracks	num_of_sales	rolling_stone_critic	mtv_critic	music_r	maniac_crit
1	1767	Call me Cat	Moneyless That Doggies		Folk	2	2006	11		905193	4	1.5	3	
2	23548	Down Mare			Metal	2	014	7		969122	3	4	5	
3	17822	Embarrasse	d Hungry		Latino	2	2000	11		522095	2.5	1	2	
4	19565	Standard Im	mediate Ennineer Slovakia		Pon	2	n17	4		610116	15	2	4	
id	real_nam	ne	art_name	role		year_of_	birth	country	cty		email			zip_code
1	Kameko I	Nelson		female	voice	1981		Slovenia	Vedrin		lorem.tristique.aliquet@nonummyFus		ce.co.uk	6914
2	Sacha Ca	ameron	Bianco Ibureno Chloride	artist		1955	4	Guernsey	Montig	ny-le-Tilleul	Sed@elementum.com	n		0568
3	Thane E.	Oliver		rapper		1993	4	Guinea-Bissau	Saint-S	ébastien-sur-Loire	Integer.tincidunt.aliqu	am@libero.com		813154
4	Cole E. J	oseph	Persian responsible	bassis	t	1994		Estonia	Newau	av	nulla.magna.malesua	da@vulputate.co	om	A4S 2B3

Expected format of results:

Erin F. Evans Norway 5.0

Hayley Russell Norway 4.8

Ahmed Rosa Norway 3.2

Explore using Spark SQL and Dataset API.

id	artist_id	album_title	genre	year_of_pub	num_of_tracks	num_of_sales	rolling_stone_critic	mtv_critic	music_maniac_critic
1	1767	Call me Cat Moneyless That Doggies	Folk	2006	11	905193	4	1.5	3
2	23548	Down Mare	Metal	2014	7	969122	3	4	5
3	17822	Embarrassed Hungry	Latino	2000	11	522095	2.5	1	2
4	19565	Standard Immediate Engineer Slovakia	Pon	2017	4	610116	15	2	4

id	real_name	art_name	role	year_of_birth	country	city	email	zip_code
1	Kameko Nelson		female voice	1981	Slovenia	Vedrin	lorem.tristique.aliquet@nonummyFusce.co.uk	6914
2	Sacha Cameron	Bianco Ibureno Chloride	artist	1955	Guernsey	Montigny-le-Tilleul	Sed@elementum.com	0568
3	Thane E. Oliver		rapper	1993	Guinea-Bissau	Saint-Sébastien-sur-Loire	Integer.tincidunt.aliquam@libero.com	813154
4	Cole E. Joseph	Persian responsible	bassist	1994	Estonia	Newauav	nulla.magna.malesuada@vulputate.com	A4S 2B3

Expected format of results:

450

1589

QUESTIONS?