Hive作业

辅助资料

- · 基本SQL教程:
- https://www.w3school.com.cn/sql/index.asp
- https://www.liaoxuefeng.com/wiki/1177760294764384

- Hive DDL:
- https://cwiki.apache.org/confluence/display/Hive/LanguageManual +DDL

数据解释-hive_sql_test1.t_user表

- t_user观众表共6000+条数据
- 字段为: UserID, Sex, Age, Occupation, Zipcode
- 字段中文解释:用户id,性别,年龄,职业,邮编

		t_user.userid	t_user.sex	t_user.age	t_user.occupation	t_user.zipcode
ili v	1	1	F	1	10	48067
	2	2	М	56	16	70072
1	3	3	М	25	15	55117
	4	4	М	45	7	02460
	5	5	М	25	20	55455
	6	6	F	50	9	55117
	7	7	М	35	1	06810
	8	8	М	25	12	11413
	9	9	М	25	17	61614
	10	10	F	35	1	95370

数据解释-hive_sql_test1.t_movie表

- t_movie电影表共3000+条数据
- 字段为: MovieID, MovieName, MovieType
- 字段中文解释: 电影ID, 电影名, 电影类型

		t_movie.movieid	t_movie.moviename	t_movie.movietype
	1	1	Toy Story (1995)	Animation Children's Comedy
	2	2	Jumanji (1995)	Adventure Children's Fantasy
±	3	3	Grumpier Old Men (1995)	Comedy Romance
	4	4	Waiting to Exhale (1995)	Comedy Drama
	5	5	Father of the Bride Part II (1995)	Comedy
	6	6	Heat (1995)	Action Crime Thriller
	7	7	Sabrina (1995)	Comedy Romance
	8	8	Tom and Huck (1995)	Adventure Children's
	9	9	Sudden Death (1995)	Action
	10	10	GoldenEye (1995)	Action Adventure Thriller

数据解释-hive_sql_test1.t_rating表

- t_rating影评表100万+条数据
- 字段为: UserID, MovieID, Rate, Times
- 字段中文解释:用户ID,电影ID,评分,评分时间

		t_rating.userid	t_rating.movieid	t_rating.rate	t_rating.times
	1	1	1193	5	978300760
	2	1	661	3	978302109
土	3	1	914	3	978301968
	4	1	3408	4	978300275
	5	1	2355	5	978824291
	6	1	1197	3	978302268
	7	1	1287	5	978302039
	8	1	2804	5	978300719
	9	1	594	4	978302268
	10	1	919	4	978301368

题目一

简单:展示电影ID为2116这部电影各年龄段的平均影评分

INFO : Completed executing command(queryId=hive_20210802153200_beb2a6f3-7f1f-49c1-9d86-584ff07b976b); Time taken: 39.145 seconds

INFO : OK

查询	同历史 [·]	记录	保存的查询	结果 (7)
		age		avgrate
ılıl 🔻	1	1		3.2941176470588234
	2	18		3.3580246913580245
£	3	25		3.436548223350254
	4	35		3.2278481012658227
	5	45		2.8275862068965516
	6	50		3.32
	7	56		3.5

题目二

中等:找出男性评分最高且评分次数超过50次的10部电影,展示电影名,平均影评分和评分次数

	sex	name	avgrate	total
1	М	Sanjuro (1962)	4.639344262295082	61
2	М	Godfather, The (1972)	4.583333333333333	1740
3	М	Seven Samurai (The Magnificent Seven) (Shichinin no samurai) (1954)	4.576628352490421	522
4	М	Shawshank Redemption, The (1994)	4.560625	1600
5	М	Raiders of the Lost Ark (1981)	4.520597322348094	1942
6	М	Usual Suspects, The (1995)	4.518248175182482	1370
7	М	Star Wars: Episode IV - A New Hope (1977)	4.495307167235495	2344
8	М	Schindler's List (1993)	4.49141503848431	1689
9	М	Paths of Glory (1957)	4.485148514851486	202
10	М	Wrong Trousers, The (1993)	4.478260869565218	644

题目三(选做)

困难:找出影评次数最多的女士所给出最高分的10部电影的平均影评分,展示电影名和平均影评分(可使用多行SQL)

t. moviename	t. avgrate
Close Shave, A (1995)	4. 52054794520548
Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb (1963)	4. 4498902706656913
Rear Window (1954)	4. 476190476190476
It Happened One Night (1934)	4. 280748663101604
Crying Game, The (1992)	3.7314890154597236
Trust (1990)	4. 188888888888888
Duck Soup (1933)	4. 2104377104377
Night on Earth (1991)	3.747422680412371
Roger & Me (1989)	4.073934837092731
Being John Malkovich (1999)	4.125390450691656

作业提交

在HUE上使用student账户跑SQL

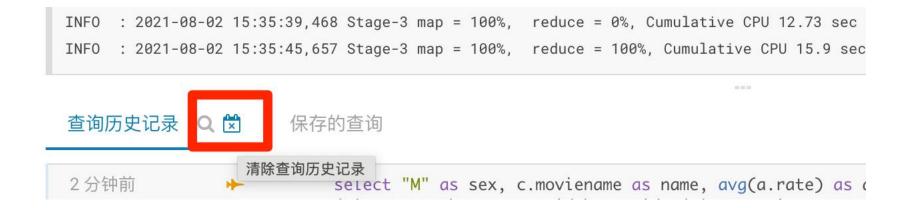
截图上缴作业,要求SQL,如图queryld和结果在一张截图里

(PS: 下图的sql和结果没有截全)



其他

在HUE上只有student一个账户,大家记得跑完清除历史记录!



提交链接及截止时间

- 作业提交链接: https://jinshuju.net/f/DWOaL1
- 作业截止日期: 8 月 8 日 23:59 前

附加作业: GeekFileFormat

请优先完成前面三个作业,Hive的练习更多的是对于HQL的使用,对于完成的同学,可以试着写一个Hive的FileFormat: GeekFileFormat

要求:

- 实现两个类: GeekTextInputFormat和GeekTextOutputFormat
- 建表时使用create table ... stored as geek来创建GeekFormat表
- 该表的文件类型为文本类型,非二进制类型
- 类似Base64TextInputFormat和Base64TextOutputFormat,GeekFormat也是用于加密
- 解密规则如下:文件中出现任何的geek,geeeek,geeeeeeeeeek等**单词**时,进行过滤,即删除该单词。gek需要保留。字母中连续的"e"最大长度为256个。
 - 例如: This notebook can be geeeek used to geek install gek on all geeeek worker nodes, run data generation, and create the TPCDS geeeeeeeek database.
 - 解密为: This notebook can be used to install gek on all worker nodes, run data generation, and create the TPCDS database.
- 【附加的附加】加密规则如下:文件输出时每随机2到256个单词,就插入一个gee...k,字母e的个数等于前面出现的非gee...k单词的个数。
 - 例如: This notebook can be used to install gek on all worker nodes, run data generation, and create the TPCDS database.
 - 加密为: This notebook can be geeeek used to geek install gek on all geeeek worker nodes, run data generation, and create the TPCDS geeeeeeeek database.