

# Solution – Exercise II

HiveQL, Create and work with External Tables on  
IMDb Data



# Solution

## Prerequisites:

- Setup Google Cloud SDK
- Start VM instance
- Pull docker container `marcelmittelstaedt/hive_base:latest`
- Start docker container: `docker run -dit --name hive_base_container -p 8088:8088 -p 9870:9870 -p 9864:9864 marcelmittelstaedt/hive_base:latest`
- Get into docker container
- Start Hadoop and Hive Shell:
  - `start-all.sh`
  - `hive`

# Solution

## Exercise 1-4:

1. Download and unzip <https://datasets.imdbws.com/name.basics.tsv.gz>

```
wget https://datasets.imdbws.com/name.basics.tsv.gz
gunzip name.basics.tsv.gz
```

2. Create HDFS directory **/user/hadoop/imdb/name\_basics/** for file name.basics.tsv

```
hadoop fs -mkdir /user/hadoop/imdb/name_basics
```

3. Put TSV file to HDFS:

```
hadoop fs -put name.basics.tsv /user/hadoop/imdb/name_basics/name.basics.tsv
```

# Solution

## Exercise 1-4:

### 4. Create Hive Table `name_basics`:

```
hive > CREATE EXTERNAL TABLE IF NOT EXISTS name_basics(  
    nconst STRING,  
    primary_name STRING,  
    birth_year INT,  
    death_year STRING,  
    primary_profession STRING,  
    known_for_titles STRING  
    ) COMMENT 'IMDb Actors' ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' ST  
ORED AS TEXTFILE LOCATION '/user/hadoop/imdb/name_basics'  
TBLPROPERTIES ('skip.header.line.count'='1');
```

# Solution

## Exercise 5:

*a) How many movies and how many TV series are within the IMDB dataset?*

```
hive > SELECT m.title_type, count(*)  
       FROM title_basics m GROUP BY m.title_type;  
  
tvMovie 137831  
movie 623469  
tvEpisode 7002447  
tvSeries 232286  
[...]  
  
Time taken: 32.908 seconds, Fetched: 11 row(s)
```

*b) Who is the youngest actor/writer/... within the dataset?*

```
hive > SELECT * FROM name_basics n  
       WHERE n.birth_year = ( SELECT MAX(birth_year) FROM name_basics);
```

# Solution

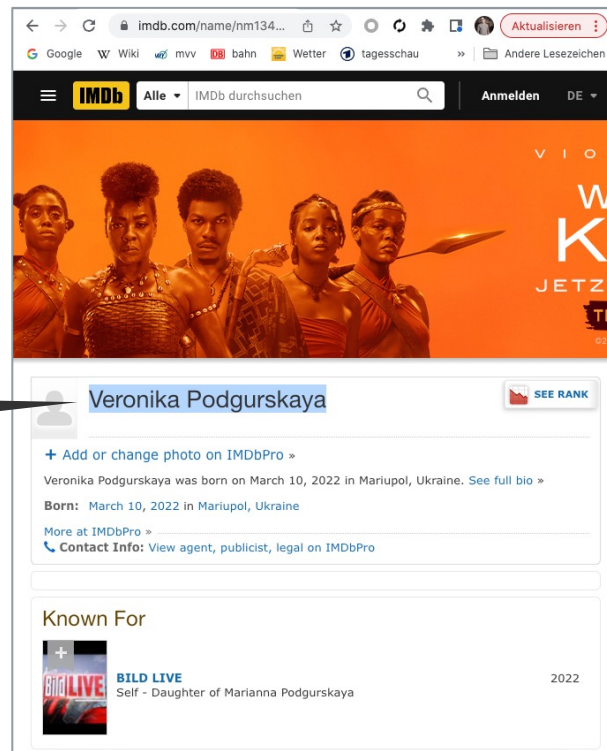
## Exercise 5:

*b) Who is the youngest actor/writer/... within the dataset?*

```
hive > SELECT * FROM name_basics n
      WHERE n.birth_year = ( SELECT MAX(birth_year)
                             FROM name_basics);
```

And it's **Veronika Podgurskaya**, daughter of Marianna Podgurskaya, born March 10, 2022 in Mariupol and kidnapped by russians during russian illegal war of aggression against **Ukraine**.

```
nm13478983 Veronika Podgurskaya 2022 NULL NULL
nm13514183 Jacques Webster 2022 NULL NULL
nm13607196 Kira Glodan 2022 2022 NULL
nm13810098 Zane Green 2022 NULL NULL
nm13945207 Fritz 2022 NULL NULL
Time taken: 65.166 seconds, Fetched: 5 row(s)
```



# Solution

## Exercise 5:

- c) Create a list (*m.tconst*, *m.original\_title*, *m.start\_year*, *r.average\_rating*, *r.num\_votes*) of movies which are:
- equal or newer than year 2010
  - have an average rating equal or better than 8,1
  - have been voted more than 100.000 times

```
hive > SELECT m.tconst, m.original_title, m.start_year, r.average_rating, r.num_votes
FROM title_basics m JOIN title_ratings r on (m.tconst = r.tconst)
WHERE r.average_rating >= 8.1 and m.start_year >= 2010 and m.title_type = 'movie'
and r.num_votes > 100000
ORDER BY r.average_rating desc, r.num_votes DESC;
```

```
tt15097216 Jai Bhim 2021 8.9 198911
tt1375666 Inception 2010 8.8 2321825
tt10189514 Soorarai Pottru 2020 8.7 114950
tt0816692 Interstellar 2014 8.6 1791616
tt1675434 Intouchables 2011 8.5 848765
tt2582802 Whiplash 2014 8.5 842994
tt6751668 Gisaengchung 2019 8.5 780914
tt1345836 The Dark Knight Rises 2012 8.4 1687333
tt1853728 Django Unchained 2012 8.4 1532742
tt7286456 Joker 2019 8.4 1249071
[...]
```

# Solution

## Exercise 5:

*d) How many movies are in list of c)?*

```
hive > SELECT count(*)  
        FROM title_basics m JOIN title_ratings r on (m.tconst = r.tconst)  
        WHERE r.average_rating >= 8.1 and m.start_year >= 2010 and m.title_type = 'movie'  
        and r.num_votes > 100000;
```

55



# Solution

## Exercise 5:

e) *We want to know which years have been great for cinema.*

*Create a list with one row per year and a related count of movies which:*

- have an average rating better than 8*
  - have been voted more than 100.000 times*
- ordered descending by count of movies.*

```
hive > SELECT m.start_year, count(*)  
        FROM title_basics m JOIN title_ratings r on (m.tconst = r.tconst)  
        WHERE r.average_rating > 8 AND m.title_type = 'movie'  
        AND r.num_votes > 100000  
        GROUP BY m.start_year  
        ORDER BY count(*) DESC;
```

```
1995 8  
1994 6  
2019 6  
2009 6  
2014 6  
2001 6  
[...]
```

# Solution

## Exercise 5:

*So 1995 seems to be a really good year for cinema, 8 really good movies have been releases, but which are they?*

```
hive > SELECT
        m.tconst, m.original_title, m.start_year, r.average_rating,
        r.num_votes
FROM title_basics m JOIN title_ratings r ON (m.tconst = r.tconst)
WHERE
        r.average_rating > 8 AND m.title_type = 'movie'
        AND r.num_votes > 100000 AND m.start_year = 1995
ORDER BY r.average_rating DESC;

tt0114369 Se7en 1995 8.6 1629859
tt0114814 The Usual Suspects 1995 8.5 1076821
tt0112573 Braveheart 1995 8.4 1031223
tt0114709 Toy Story 1995 8.3 984966
tt0113277 Heat 1995 8.3 650043
tt0112641 Casino 1995 8.2 516258
tt0113247 La haine 1995 8.1 172841
tt0112471 Before Sunrise 1995 8.1 306176

[...]
```