## Spark Preparation

We check if we are in Google Colab. If this is the case, install all necessary packages.

To run spark in Colab, we need to first install all the dependencies in Colab environment i.e. Apache Spark 3.3.2 with hadoop 3.3, Java 8 and Findspark to locate the spark in the system. The tools installation can be carried out inside the Jupyter Notebook of the Colab. Learn more from [A Must-Read Guide on How to Work with PySpark on Google Colab for Data Scientists!](#)

```
1 try:
2   import google.colab
3   IN_COLAB = True
4 except:
5   IN_COLAB = False
```

```
1 if IN_COLAB:
2     !apt-get install openjdk-8-jdk-headless -qq > /dev/null
3     !wget -q https://dlcdn.apache.org/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz
4     !tar xf spark-3.5.1-bin-hadoop3.tgz
5     !mv spark-3.5.1-bin-hadoop3 spark
6     !pip install -q findspark
7     import os
8     os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
9     os.environ["SPARK_HOME"] = "/content/spark"
```

## Start a Local Cluster

```
1 import findspark
2 findspark.init()
```

```
1 spark_url = 'local'
```

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql import SQLContext
3
4 spark = SparkSession.builder\
5         .master(spark_url)\
6         .appName('Spark Assignment')\
7         .config('spark.ui.port', '4040')\
8         .getOrCreate()
```

## Spark Assignment

Based on the movie review dataset in 'netflix-rotten-tomatoes-metacritic-imdb.csv', answer the below questions.

**Note:** do not clean or remove missing data

```
1 sc = spark.sparkContext
```

```
1 spark
```

**SparkSession - in-memory**

**SparkContext**

[Spark UI](#)

Version
      v3.5.1
Master
      local
AppName
      Spark Assignment

```
1 sc
```

**SparkContext**
[Spark UI](#)

Version
        v3.5.1
Master
        local
AppName
        Spark Assignment

```
1 data_path = './netflix-rotten-tomatoes-metacritic-imdb.csv'
```

```
1 df = spark.read.csv(data_path, header=True, inferSchema=True)
```

## ⌄ What is the maximum and average of the overall hidden gem score?

```
1 from pyspark.sql.functions import max, avg
2
3 hidden_gem_stats = df.agg(
4     max("Hidden Gem Score").alias("Maximum Hidden Gem Score"),
5     avg("Hidden Gem Score").alias("Average Hidden Gem Score")
6 )
7
8 hidden_gem_stats.show()
```

```
+------------------------+------------------------+
|Maximum Hidden Gem Score|Average Hidden Gem Score|
+------------------------+------------------------+
|                     9.8|       5.937551386501226|
+------------------------+------------------------+
```

## ⌄ How many movies that are available in Korea?

```
1 from pyspark.sql.functions import col
2
3 korea_movie_count = df.filter(col("Languages").contains("Korea")).count()
4 print(f"Number of movies available in Korea: {korea_movie_count}")
```

```
Number of movies available in Korea: 735
```

## ⌄ Which director has the highest average hidden gem score?

```
1 from pyspark.sql.functions import avg, col
2
3 director_avg_scores = df.groupBy("Director") \
4                         .agg(avg("Hidden Gem Score").alias("Average Hidden Gem Score")) \
5                         .orderBy(col("Average Hidden Gem Score").desc())
6
7 top_director = director_avg_scores.first()
8
9 print(f"Director with the highest average hidden gem score: {top_director['Director']} (Score: {top_director['Average Hi
```

```
Director with the highest average hidden gem score: Dorin Marcu (Score: 9.8)
```

## ⌄ How many genres are there in the dataset?

```
1 from pyspark.sql.functions import explode, split, col
2
3 genre_df = df.withColumn("Genre", explode(split(col("Genre"), ",\s*")))
4 genre_count = genre_df.select("Genre").distinct().count()
5
6 print(f"Total number of unique genres: {genre_count}")
```

```
Total number of unique genres: 28
```