# 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-3.3.2/spark-3.3.2-bin-hadoop3.tgz
4     !tar xf spark-3.3.2-bin-hadoop3.tgz
5     !mv spark-3.3.2-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.functions import col
```

```
1 spark = SparkSession.builder\
2         .master(spark_url)\
3         .appName('Spark SQL')\
4         .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 path = 'netflix-rotten-tomatoes-metacritic-imdb.csv'
```

```
1 df = spark.read.option("header", True).csv(path)
```

```
1 df.show(3)
```

```
+-----------------+-----------------+-----------------+---------------+--------------+----------------+-
|            Title|            Genre|             Tags|      Languages|Series or Movie|Hidden Gem Score|
+-----------------+-----------------+-----------------+---------------+--------------+----------------+-
|  Lets Fight Ghost|Crime, Drama, Fan...|Comedy Programmes...|Swedish, Spanish|         Series|             4.3|
|HOW TO BUILD A GIRL|          Comedy|Dramas,Comedies,F...|        English|          Movie|             7.0|
|        Centigrade|  Drama, Thriller|        Thrillers|        English|          Movie|             6.4|
+-----------------+-----------------+-----------------+---------------+--------------+----------------+-
only showing top 3 rows
```

```
1 df.columns
```

```
['Title',
 'Genre',
```

```
'Tags',
'Languages',
'Series or Movie',
'Hidden Gem Score',
'Country Availability',
'Runtime',
'Director',
'Writer',
'Actors',
'View Rating',
'IMDb Score',
'Rotten Tomatoes Score',
'Metacritic Score',
'Awards Received',
'Awards Nominated For',
'Boxoffice',
'Release Date',
'Netflix Release Date',
'Production House',
'Netflix Link',
'IMDb Link',
'Summary',
'IMDb Votes',
'Image',
'Poster',
'TMDb Trailer',
'Trailer Site']
```

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

```
1 df = df.withColumn('Hidden Gem Score', df['Hidden Gem Score'].cast('float'))
```

```
1 from pyspark.sql.functions import avg, min, max, countDistinct
2 df.select(avg('Hidden Gem Score'), max('Hidden Gem Score')).show()
```

```
+--------------------+--------------------+
|avg(Hidden Gem Score)|max(Hidden Gem Score)|
+--------------------+--------------------+
|    5.937551392466598|                 9.8|
+--------------------+--------------------+
```

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

### Consider "Language" column

```
1 from pyspark.sql.functions import array_contains
2 from pyspark.sql.functions import split
3 from pyspark.sql.functions import split, explode
4
5 movies_df = df.withColumn("Language", split("Languages", ",\s*"))
6 korean_movies_df = movies_df.filter(array_contains("Language", "Korean"))
7 # korean_movies_df = korean_movies_df.filter(col("Series or Movie") == "Movie")
8 korean_movie_count = korean_movies_df.distinct().count()
9 print("Number of movies available in Korea:", korean_movie_count)
10
```

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

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

```
1 director_scores_df = df.groupBy("Director").agg(avg("Hidden Gem Score").alias("Avg Hidden Gem Score"))
2 sorted_director_scores_df = director_scores_df.sort(col("Avg Hidden Gem Score").desc()).first()
3 top_director = sorted_director_scores_df["Director"]
4 top_score = sorted_director_scores_df["Avg Hidden Gem Score"]
5 print("The director with the highest average hidden gem score is", top_director, "with an average score of", top_sco
6
```

```
The director with the highest average hidden gem score is Dorin Marcu with an average score of 9.800000190734863
```

```
1 dir_df = df.withColumn("Director", split(df["Director"], ",\s*"))
2 dir_df = dir_df.select(explode(dir_df["Director"]).alias("Director"), dir_df["Hidden Gem Score"])
3 director_scores = dir_df.groupBy("Director").agg({"Hidden Gem Score": "avg"})
4 top_director = director_scores.orderBy(director_scores["avg(Hidden Gem Score)"].desc()).first()
5 print(f"The director with the highest average Hidden Gem Score is {top_director['Director']} with an average score o
```

    The director with the highest average Hidden Gem Score is Dorin Marcu with an average score of 9.800000190734863

## How many genres are there in the dataset?

```
1 genres_df = df.select(split(col("Genre"), ",\s*").alias("Genres_array"))
2 exploded_genres_df = genres_df.select(explode(col("Genres_array")).alias("Genres"))
3 num_genres = exploded_genres_df.select("Genres").distinct().count()
4 print("There are", num_genres, "genres in the dataset.")
```

    There are 28 genres in the dataset.

✓  0s    completed at 8:35 PM

```
1 genres_df = df.select(split(col("Genre"), ",\s*").alias("Genres_array"))
```