```
Aluna: Samantha Pinheiro Gomes
 RM: 362698
 Turma: 29ABD
 Matéria: Distributed Data Processing & Storage
 Este notebook realiza a análise do arquivo "results", que contém dados de diversos jogos realizados em várias cidades e países ao redor do mundo, incluindo as datas dos eventos e seus respectivos placares.
 Import de bibliotecas
 from pyspark.sql import SparkSession
   import pandas as pd
   import requests
   from pyspark.sql.types import StructType, StructField, StringType, IntegerType, DoubleType, DateType, BooleanType
   from io import StringIO
   from pyspark.sql.functions import to_date, col, max
 ▶ √ 07:02 PM (<1s)
   spark =SparkSession.builder \
        .appName("Analise arquivo results") \
        .getOrCreate()
 Baixa arquivo armazenado no GitHub
   url = "https://raw.githubusercontent.com/samppinheiro/fiap_distributed_data_processing_-_storage/main/results.csv"
   response = requests.get(url)
   csv_content = response.content.decode("utf-8")
   pandas_df = pd.read_csv(StringIO(csv_content))
▶ ■ pandas_df: pandas.core.frame.DataFrame = [date: object, home_teamName: object ... 7 more fields]
 Load do dados em um DF do spark
    schema = StructType([
       StructField("date", StringType(), True),
       StructField("home_team", StringType(), True),
       StructField("away_team", StringType(), True),
       StructField("home_score", IntegerType(), True),
       StructField("away_score", IntegerType(), True),
       StructField("tournament", StringType(), True),
       StructField("city", StringType(), True),
       StructField("country", StringType(), True),
        StructField("neutral", BooleanType(), True)
   1)
   df_results = spark.createDataFrame(pandas_df, schema=schema)
    df_results = df_results.withColumn("date", to_date("date", "yyyy-MM-dd"))
▶ ■ df_results: pyspark.sql.connect.dataframe.DataFrame = [date: date, home_team: string ... 7 more fields]
 Define dataframe para iniciar a análise
   df_results_consol = (
        df_results
        .withColumnRenamed("date", "dt_torneio")
        .withColumnRenamed("home_team", "nm_time_casa")
        .withColumnRenamed("away_team", "nm_time_visitante")
        .withColumnRenamed("home_score", "nr_placar_casa")
        .withColumnRenamed("away_score", "nr_placar_visitante")
        .withColumnRenamed("tournament", "nm_torneio")
        .withColumnRenamed("city", "nm_cidade")
        .withColumnRenamed("country", "nm_pais")
   ).drop("neutral")
 ▶ ■ df_results_consol: pyspark.sql.connect.dataframe.DataFrame = [dt_torneio: date, nm_time_casa: string ... 6 more fields]
 Análises realizadas no arquivo
    nm_registro = df_results_consol.count()
   display(f"Existem {nm_registro} registros na base")
                                                                                                                                                                                                                                            Optimize
'Existem 40839 registros na base'
   ## Quantas equipes únicas mandantes existem na base?
   nm_equipes_mandantes = df_results_consol.select("nm_time_casa").distinct().count()
   display(f"Existem {nm_equipes_mandantes} equipes mandantes na base")
 > <u>Illi</u> See performance (1)
                                                                                                                                                                                                                                            Optimize
'Existem 309 equipes mandantes na base'
   nm_vitorias_mandantes = (df_results_consol
                            .filter(col("nr_placar_casa") > col("nr_placar_visitante"))
                            .count())
   display(f"As equipes mandantes saíram vitoriosas {nm_vitorias_mandantes} vezes")
                                                                                                                                                                                                                                            Optimize
'As equipes mandantes saíram vitoriosas 19864 vezes'
   nm_vitorias_visitantes = (df_results_consol
                            .filter(col("nr_placar_visitante") > col("nr_placar_casa"))
                            .count())
   display(f"As equipes visitantes saíram vitoriosas {nm_vitorias_visitantes} vezes")
                                                                                                                                                                                                                                            Optimize
'As equipes visitantes saíram vitoriosas 11544 vezes'
   ## Quantas partidas resultaram em empate?
    nm_empate = (df_results_consol
                            .filter(col("nr_placar_casa") == col("nr_placar_visitante"))
                            .count())
   display(f'{nm_empate} partidas resultaram em empate')
                                                                                                                                                                                                                                            Optimize
'9431 partidas resultaram em empate'
    df_pais = (df_results_consol
               .groupBy("nm_pais")
               .orderBy(col("count").desc())
               .withColumnRenamed("count", "qt_partidas"))
   df_pais.show()
 > <u>Illi</u> See performance (1)
                                                                                                                                                                                                                                            Optimize
 ▶ ■ df_pais: pyspark.sql.connect.dataframe.DataFrame = [nm_pais: string, qt_partidas: long]
              France
                             801
             England
                            687
            Malaysia
                            644
              Sweden
                            637
                             581
             Germany
              Brazil
                            529
                            517
               Spain
             Thailand
                             483
               Italy|
                             480
          Switzerland
                             477
             Austria
                             475
|United Arab Emirates|
                             472
         South Africa
                             470
               Qatar
                             467
                            453
          South Korea
           Argentina
                             449
                             431
             Hungary
               Chile
                             405
             Belgium
                             396
only showing top 20 rows
   df_partidas = df_pais.orderBy(col("qt_partidas").desc()).limit(1)
   df_partidas.show()
 > <u>Illi</u> See performance (1)
                                                                                                                                                                                                                                            Optimize
▶ ■ df_partidas: pyspark.sql.connect.dataframe.DataFrame = [nm_pais: string, qt_partidas: long]
+-----+
     nm_pais|qt_partidas|
+-----
|United States| 1144|
+-----+
    max_nr_gols = (df_results_consol
                   .withColumn("nr_maior_placar", col("nr_placar_casa") + col("nr_placar_visitante"))
```

```
.orderBy(col("nr_maior_placar").desc())
              .select("dt_torneio", "nm_time_casa", "nm_time_visitante", "nr_maior_placar")
  max_nr_gols.show()
 > <u>Illi</u> See performance (1)
                                                                                                                                                                                     Optimize
max_nr_gols: pyspark.sql.connect.dataframe.DataFrame = [dt_torneio: date, nm_time_casa: string ... 2 more fields]
+-----
|dt_torneio|nm_time_casa|nm_time_visitante|nr_maior_placar|
+----+
|2001-04-11| Australia| American Samoa|
+-----+
```

```
max_diff_gols = (df_results_consol
            .withColumn("nr_maior_diff_gols", col("nr_placar_casa") - col("nr_placar_visitante"))
           .orderBy(col("nr_maior_diff_gols").desc())
           .select('*')
           .limit(1)
  max_diff_gols.show()
> <u>Illi</u> See performance (1)
                                                                                                                                                 Optimize
▶ ■ max_diff_gols: pyspark.sql.connect.dataframe.DataFrame = [dt_torneio: date, nm_time_casa: string ... 7 more fields]
|dt_torneio|nm_time_casa|nm_time_visitante|nr_placar_casa|nr_placar_visitante|
                                                    nm torneio | nm cidade | nm pais | nr maior diff gols |
|2001-04-11| Australia| American Samoa|
                                             0|FIFA World Cup qu...|Coffs Harbour|Australia|
+------
```

```
## Quantos jogos ocorreram no Brasil?
   df_brasil = (df_results_consol
                .filter(col("nm_pais") == "Brazil")
                .count()
   display(f"{df_brasil} jogos foram realizados no Brasil")
> Illi See performance (1)
'529 jogos foram realizados no Brasil'
                                                                                                        + Code → Text → Assistant
```