Bundesliga Soccer Analysis using Pyspark **Questions to Answer** Who are the winners of the D1 division in the Germany Football Association (Bundesliga) in the last decade? Which teams have been relegated in the past 10 years? Does Octoberfest affect the performance of Bundesliga? #importing required Libraries import findspark from pyspark import SparkContext from pyspark.sql import SparkSession, Window, Row from pyspark.sql.functions import * from pyspark.sql.types import * import matplotlib.pyplot as plt #https://spark.apache.org/docs/latest/sql-getting-started.html spark = SparkSession \ .builder \ .appName("firstSpark") \ .getOrCreate() # df = spark.read.format('csv').options(header='true').load(".") def load dataframe(filename): df = spark.read.format('csv').options(header='true').load(filename) return df In [4]: #creating a dataframe df matches = load dataframe('/Users/ruturajharal/Desktop/Bundesliga-analysis-pyspark/Matches.csv') df matches.limit(10).show() +-----|Match ID|Div|Season| Date| HomeTeam| AwayTeam|FTHG|FTAG|FTR| 1| D2| 2009|2010-04-04| Oberhausen|Kaiserslautern| 2| 1| H| 2| D2| 2009|2009-11-01| Munich 1860|Kaiserslautern| 0| 1| A| 3| D2| 2009|2009-10-04| Frankfurt FSV|Kaiserslautern| 1| 1| D| 4| D2| 2009|2010-02-21| Frankfurt FSV| Karlsruhe| 2| 1| H| 5| D2| 2009|2009-12-06| Ahlen| Karlsruhe| 1| 3| A| 6| D2| 2009|2010-04-03| Union Berlin| Karlsruhe| 1| 1| D| 7| D2| 2009|2009-08-14| Paderborn| Karlsruhe| 2| 0| H| 8| D2| 2009|2010-03-08| Bielefeld| Karlsruhe| 0| 1| A| 9| D2| 2009|2009-09-26|Kaiserslautern| Karlsruhe| 2| 0| H| 10| D2| 2009|2009-11-21| Hansa Rostock| Karlsruhe| 2| 1| H| #converting to pandas dataframe df_matches.limit(10).toPandas() Match_ID Div Season Date HomeTeam AwayTeam FTHG FTAG FTR 1 D2 2009 2010-04-04 Oberhausen Kaiserslautern Η 2009 2009-11-01 2 D2 Munich 1860 Kaiserslautern 3 D2 2009 2009-10-04 Frankfurt FSV Kaiserslautern 1 D 4 D2 2009 2010-02-21 Frankfurt FSV Karlsruhe 5 D2 2009 2009-12-06 Ahlen Karlsruhe 3 2009 2010-04-03 6 D2 Union Berlin Karlsruhe D2 2009 2009-08-14 Paderborn Karlsruhe 0 Н 2009 2010-03-08 Bielefeld Karlsruhe 8 Kaiserslautern 2 9 D2 2009 2009-09-26 Karlsruhe 0 Н 10 2009 D2 2009-11-21 Hansa Rostock Karlsruhe Н df matches.printSchema() root |-- Match ID: string (nullable = true) |-- Div: string (nullable = true) |-- Season: string (nullable = true) |-- Date: string (nullable = true) |-- HomeTeam: string (nullable = true) |-- AwayTeam: string (nullable = true) |-- FTHG: string (nullable = true) |-- FTAG: string (nullable = true) |-- FTR: string (nullable = true) #Next step would be renaming some of the columns old cols = df matches.columns[-3:] new_cols = ["HomeTeamGoals", "AwayTeamGoals", "FinalResult"] #combine 2 rdd's old_new_cols = [*zip(old_cols, new_cols)] for old_col, new_col in old_new_cols: df_matches = df_matches.withColumnRenamed(old_col, new_col) #checking the pandas structure df matches.limit(5).toPandas() Match_ID Div Season AwayTeam HomeTeamGoals AwayTeamGoals FinalResult Date **HomeTeam** 0 Н D2 2009 2010-04-04 Oberhausen Kaiserslautern 2 1 1 1 2009 2009-11-01 0 D2 Munich 1860 Kaiserslautern Α 3 D2 2009 2009-10-04 Frankfurt FSV Kaiserslautern 1 D 3 D2 2009 2010-02-21 Frankfurt FSV Karlsruhe Н 2009 2009-12-06 3 D2 Ahlen Karlsruhe Α Who are the winners Bundesliga in the last decade? df matches.limit(10).toPandas() **HomeTeam** Match_ID Div **Date** AwayTeam HomeTeamGoals AwayTeamGoals FinalResult Season D2 2009 2010-04-04 Oberhausen Kaiserslautern 2 Н Kaiserslautern D2 2009 2009-11-01 Munich 1860 0 2 3 1 D D2 2009 2009-10-04 Frankfurt FSV Kaiserslautern 1 3 D2 2009 2010-02-21 Frankfurt FSV 4 Karlsruhe Н 2009-12-06 4 2009 Ahlen 1 3 Karlsruhe Α D2 2009 2010-04-03 Union Berlin Karlsruhe D 6 2009 2 0 7 D2 2009-08-14 Paderborn Karlsruhe Н 2009 Bielefeld 0 8 D2 2010-03-08 Karlsruhe Α 2009-09-26 2 D2 2009 Kaiserslautern Karlsruhe Н Karlsruhe 10 D2 2009 2009-11-21 Hansa Rostock Н The approach is to aggregate the home and away game results separately creating two dataframes: home and away. #If Home team Win then result = 1, Away team win result = 1, Draw result =1 df matches = df matches.withColumn('HomeTeamWin', when(col('FinalResult') == 'H', 1).otherwise(0)) \ .withColumn('AwayTeamWin', when(col('FinalResult') == 'A', 1).otherwise(0)) \ .withColumn('GameTie', when(col('FinalResult') == 'D', 1).otherwise(0)) df matches.limit(10).toPandas() Match_ID Div Season Date **HomeTeam** AwayTeam HomeTeamGoals AwayTeamGoals FinalResult HomeTeamWin AwayTea 2010-0 D2 2009 04-2 Н 1 1 Oberhausen Kaiserslautern 04 2009-2 D2 2009 Munich 1860 Kaiserslautern 0 0 11-01 2009-0 2 3 D2 2009 Frankfurt FSV Kaiserslautern 1 1 D 10-04 2010-2009 Frankfurt FSV 3 D2 Karlsruhe 2 02-21 2009-4 5 D2 2009 Ahlen Karlsruhe 1 3 Α 0 12-06 2010-5 2009 Union Berlin 0 6 D2 04-Karlsruhe 1 D 03 2009-D2 2009 Paderborn Karlsruhe 2 Н 1 08-14 2010-D2 2009 0 0 7 8 Bielefeld Karlsruhe 1 Α 03-08 2009-8 D2 2009 Karlsruhe 2 Kaiserslautern Н 09-26 2009-Hansa 9 10 2009 2 D2 Karlsruhe Н 11-21 Rostock #bundesliga is a D1 division and we are interested in season <= 2010 and >= 2000 bundesliga = df matches \ .filter((col('Season') >= 2000) & (col('Season') <= 2010) & (col('Div') == 'D1')) bundesliga.limit(10).toPandas() AwayTeam HomeTeamGoals AwayTeamGoals FinalResult HomeTeamWin Match_ID Div Season Date HomeTeam **AwayTeamWi** 2010-0 21 D1 2009 Bochum Leverkusen 1 1 D 0 02-06 2009-Bayern 2009 0 22 D1 Leverkusen 1 1 D 11-22 Munich 2010-0 2 D1 2009 1 1 D 23 M'gladbach Leverkusen 05-08 2009-3 2 2 0 24 D1 2009 08-Mainz Leverkusen D 80 2009-4 2009 0 0 D 0 25 D1 Hamburg Leverkusen 10-17 2010-5 26 D1 2009 Stuttgart Leverkusen 2 Н 1 04-17 2010-Dortmund Leverkusen 03-20 2009-D1 2009 28 Schalke 04 Leverkusen 10-31 2009-8 29 D1 2009 0 0 Freiburg Leverkusen Α 08-22 2010-Werder 2009 Leverkusen 0 9 30 D1 2 D 02-21 Bremen In [14]: # home team features home = bundesliga.groupby('Season', 'HomeTeam') \ .agg(sum('HomeTeamWin').alias('TotalHomeWin'), sum('AwayTeamWin').alias('TotalHomeLoss'), sum('GameTie').alias('TotalHomeTie'), sum('HomeTeamGoals').alias('HomeScoredGoals'), sum('AwayTeamGoals').alias('HomeAgainstGoals')) \ .withColumnRenamed('HomeTeam', 'Team') home.limit(10).toPandas() Season Team TotalHomeWin TotalHomeLoss TotalHomeTie HomeScoredGoals HomeAgainstGoals 0 2005 7 Kaiserslautern 5 5 26.0 33.0 1 2006 Cottbus 6 6 5 21.0 22.0 2 2001 St Pauli 9 4 19.0 28.0 3 2005 Mainz 31.0 23.0 4 2006 9 22.0 19.0 Hamburg 4 4 Stuttgart 5 2003 29.0 13.0 9 1 6 2003 Hansa Rostock 10 6 1 34.0 18.0 7 2007 Hansa Rostock 5 8 4 17.0 21.0 8 5 2001 M'gladbach 6 6 21.0 21.0 2002 10 2 5 31.0 11.0 M'gladbach #away game features away = bundesliga.groupby('Season', 'AwayTeam') \ .agg(sum('AwayTeamWin').alias('TotalAwayWin'), sum('HomeTeamWin').alias('TotalAwayLoss'), sum('GameTie').alias('TotalAwayTie'), sum('AwayTeamGoals').alias('AwayScoredGoals'), sum('HomeTeamGoals').alias('AwayAgainstGoals')) .withColumnRenamed('AwayTeam', 'Team') away.limit(10).toPandas() Season Team TotalAwayWin TotalAwayLoss TotalAwayTie AwayScoredGoals AwayAgainstGoals 0 2005 3 10 4 21.0 38.0 Kaiserslautern 2006 Cottbus 5 17.0 27.0 2 2001 St Pauli 0 11 6 18.0 42.0 3 2005 Mainz 3 10 15.0 24.0 4 2006 Hamburg 6 21.0 18.0 6 5 2003 5 Stuttgart 5 23.0 11.0 6 2003 Hansa Rostock 2 8 7 21.0 36.0 7 2007 Hansa Rostock 3 12 13.0 31.0 8 2001 3 7 7 20.0 M'gladbach 32.0 In [18]: #season features window = ['Season'] window = Window.partitionBy(window).orderBy(col('WinPct').desc(), col('GoalDifferentials').desc()) table = home.join(away, ['Team', 'Season'], 'inner') \ .withColumn('GoalsScored', col('HomeScoredGoals') + col('AwayScoredGoals')) \ .withColumn('GoalsAgainst', col('HomeAgainstGoals') + col('AwayAgainstGoals')) \ .withColumn('GoalDifferentials', col('GoalsScored') - col('GoalsAgainst')) \ .withColumn('Win', col('TotalHomeWin') + col('TotalAwayWin')) \ .withColumn('Loss', col('TotalHomeLoss') + col('TotalAwayLoss')) \ .withColumn('Tie', col('TotalHomeTie') + col('TotalAwayTie')) \ .withColumn('WinPct', round((100* col('Win')/(col('Win') + col('Loss') + col('Tie'))), 2)) \ .drop('HomeScoredGoals', 'AwayScoredGoals', 'HomeAgainstGoals', 'AwayAgainstGoals') \ .drop('TotalHomeWin', 'TotalAwayWin', 'TotalHomeLoss', 'TotalAwayLoss', 'TotalHomeTie', 'TotalAwayTie') \ .withColumn('TeamPosition', rank().over(window)) table df = table.filter(col('TeamPosition') == 1).orderBy(asc('Season')).toPandas() table df Season GoalsScored GoalsAgainst GoalDifferentials Win Loss Tie WinPct TeamPosition Bayern Munich 2000 37.0 55.88 1 62.0 25.0 9 6 Leverkusen 2001 77.0 38.0 39.0 21 7 6 61.76 1 Bayern Munich 2002 70.0 25.0 23 5 6 67.65 1 45.0 Werder Bremen 2003 38.0 8 79.0 41.0 22 4 64.71 Bayern Munich 2004 75.0 33.0 42.0 24 5 5 70.59 1 Bayern Munich 2005 67.0 32.0 35.0 22 3 64.71 2006 7 6 Stuttgart 61.0 37.0 24.0 21 6 61.76 1 Bayern Munich 2007 21.0 22 10 64.71 7 68.0 47.0 2 Wolfsburg 8 2008 80.0 41.0 61.76 39.0 21 7 6 1 Bayern Munich 2009 72.0 31.0 41.0 20 10 58.82 1 Dortmund 10 2010 67.0 22.0 23 1 45.0 5 6 67.65 From this we can see that Bayern Munich Have won 6 times in the period of 2000-2010 which is more than 60% in 10 years Q2. Which teams were relegated in th past 2000 - 2010 In [44]: df_teams = load_dataframe('/Users/ruturajharal/Desktop/Bundesliga-analysis-pyspark/Teams.csv') df teams.limit(5).toPandas() Season TeamName KaderHome AvgAgeHome ForeignPlayersHome OverallMarketValueHome AvgMarketValueHome StadiumCapac Out[44]: Bayern 0 2017 27 26 15 597950000 22150000 75(Munich 416730000 2017 Dortmund 25 12630000 813 33 18 2 2017 Leverkusen 31 24 222600000 7180000 30: 15 180130000 3 2017 **RB** Leipzig 23 15 6000000 429 179550000 4 2017 Schalke 04 29 24 17 6190000 62 In [46]: # lets check for 2010s relegated = table.filter((col('TeamPosition') == 16) | (col('TeamPosition') == 17) | (col('TeamPosition') == 18)).orderBy(asc('Season')) relegated.filter(col('Season') == 2010).toPandas() Team Season GoalsScored GoalsAgainst GoalDifferentials Win Loss Tie WinPct TeamPosition Out[46]: Wolfsburg 43.0 48.0 2010 -5.0 26.47 16 14 1 Ein Frankfurt 2010 31.0 49.0 -18.0 18 26.47 17 St Pauli 35.0 2010 68.0 -33.0 8 21 5 23.53 18 We can see from the above table that team Wolfsburg, Ein Frankfuty and St Pauli have been relegated during season 2010