SedonaRegistrator.registerAll(spark) Ivy Default Cache set to: /root/.ivy2/cache The jars for the packages stored in: /root/.ivy2/jars :: loading settings :: url = jar:file:/home/ubuntu/spark-3.0.3-bin-hadoop3.2/jars/ivy-2.4.0.jar!/org/apache/ivy/core/settings/ivysettings.xml org.apache.sedona#sedona-python-adapter-3.0_2.12 added as a dependency org.datasyslab#geotools-wrapper added as a dependency :: resolving dependencies :: org.apache.spark#spark-submit-parent-7f20cc1f-9313-4965-bd9f-54d02e90380e;1.0 confs: [default] found org.apache.sedona#sedona-python-adapter-3.0_2.12;1.0.1-incubating in central found org.locationtech.jts#jts-core;1.18.0 in central found org.wololo#jts2geojson;0.16.1 in central found com.fasterxml.jackson.core#jackson-databind;2.12.2 in central found com.fasterxml.jackson.core#jackson-annotations;2.12.2 in central found com.fasterxml.jackson.core#jackson-core;2.12.2 in central found org.apache.sedona#sedona-core-3.0_2.12;1.0.1-incubating in central found org.apache.sedona#sedona-sql-3.0_2.12;1.0.1-incubating in central found org.datasyslab#geotools-wrapper;geotools-24.1 in central :: resolution report :: resolve 423ms :: artifacts dl 9ms :: modules in use: com.fasterxml.jackson.core#jackson-annotations;2.12.2 from central in [default] com.fasterxml.jackson.core#jackson-core;2.12.2 from central in [default] com.fasterxml.jackson.core#jackson-databind;2.12.2 from central in [default] org.apache.sedona#sedona-core-3.0_2.12;1.0.1-incubating from central in [default] org.apache.sedona#sedona-python-adapter-3.0_2.12;1.0.1-incubating from central in [default] org.apache.sedona#sedona-sql-3.0_2.12;1.0.1-incubating from central in [default] org.datasyslab#geotools-wrapper;geotools-24.1 from central in [default] org.locationtech.jts#jts-core;1.18.0 from central in [default]

Project Statement for Milestone 4

Group 6

Group member: Jinyang Ruan, Rusu Wu, Yi Yao, Brian Chan, Junqiao Mou

from pyspark.sql import SparkSession

import matplotlib.pyplot as plt

import geopandas as gpd

import seaborn as sns import numpy as np

spark = SparkSession.\ builder.\

getOrCreate()

master('local[*]').\

:: evicted modules:

confs: [default]

e.g. Find the list of Airlines having 1 stop (only showing top 10)

SELECT ' FROM routes

WHERE Stops = 1""")

There are two categories for "Stops" listed as "0" or "1"

routes_1Stop.toPandas().head(10).style

routes_1Stop = spark.sql("""

from sedona.register import SedonaRegistrator

appName('Python Spark Apache Sedona').\

from sedona.utils import SedonaKryoRegistrator, KryoSerializer

config("spark.serializer", KryoSerializer.getName).\

config("spark.kryo.registrator", SedonaKryoRegistrator.getName).\

'org.datasyslab:geotools-wrapper:geotools-24.1').\

org.wololo#jts2geojson;0.16.1 from central in [default]

______ :: retrieving :: org.apache.spark#spark-submit-parent-7f20cc1f-9313-4965-bd9f-54d02e90380e

| modules || artifacts |
conf | number| search|dwnlded|evicted|| number|dwnlded|

default | 10 | 0 | 0 | 1 || 9 | 0 |

config('spark.jars.packages','org.apache.sedona:sedona-python-adapter-3.0_2.12:1.0.1-incubating,'

org.locationtech.jts#jts-core;1.18.1 by [org.locationtech.jts#jts-core;1.18.0] in [default]

0 artifacts copied, 9 already retrieved (0kB/13ms) 21/11/25 01:27:49 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties Setting default log level to "WARN". To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel). True Out[1]: # data preparation # airports airports_dat = spark.read.option("delimiter", ",").option("header", "false").csv("airports.dat").\ toDF("AirportID", "Name", "City", "Country", "IATA", "ICAO", "Latitude", "Longitude", "Altitude", "Timezone", "DST", "Tz", "Type", "Source") airports = airports_dat.selectExpr("AirportID", "Name", "City", "Country", "IATA", "ICAO", "Latitude", "Longitude", "Altitude", "Timezone", "DST", "Tz", "Type", "Source") airports.createOrReplaceTempView("airports") # routes routes_dat = spark.read.option("delimiter", ",").option("header", "false").csv("routes.dat").\ toDF("Airline", "AirlineID", "SourceAirport", "SourceAirportID", "DestinationAirport", "DestinationAirportID", "Codeshare", "Stops", "Equipment") routes = routes_dat.selectExpr("Airline", "AirlineID", "SourceAirport", "SourceAirportID", "DestinationAirport", "DestinationAirportID", "Codeshare", "Stops", "Equipment") routes.createOrReplaceTempView("routes") # airlines airlines_dat = spark.read.option("delimiter", ",").option("header", "false").csv("airlines.dat").\ toDF("AirlineID", "Name", "Alias", "ITAT", "ICAO", "Callsign", "Country", "Active") airlines = airlines_dat.selectExpr("AirlineID", "Name", "Alias", "ITAT", "ICAO", "Callsign", "Country", "Active") airlines.createOrReplaceTempView("airlines") # states states_wkt = spark.read.option("delimiter", "\t").option("header", "false").csv("boundary-each-state.tsv").\ toDF("s_name", "s_bound") states = states_wkt.selectExpr("s_name", "ST_GeomFromWKT(s_bound) as s_bound") states.createOrReplaceTempView("states") Find the list of airports operating in the Country X e.g. Find the list of aiports operating in the United States (only showing top 10) airports_US = spark.sql(""" SELECT * FROM airports WHERE Country = 'United States'""") airports_US.toPandas().head(10).style Country IATA ICAO Latitude Longitude Altitude Timezone DST Out[3]: AirportID Name City Tz Type Source 0 3411 Barter Island LRRS Airport Barter Island United States BTI PABA 70.1340026855 -143.582000732 2 A America/Anchorage airport OurAirports 3412 35 Wainwright Air Station Fort Wainwright United States N PAWT 70.61340332 -159.8600006 A America/Anchorage airport OurAirports 3413 Cape Lisburne LRRS Airport Cape Lisburne United States LUR PALU 68.87509918 -166.1100006 16 -9 A America/Anchorage airport OurAirports 2 69.73290253 22 3414 Point Lay LRRS Airport Point Lay United States PIZ PPIZ -163.0050049 America/Anchorage airport OurAirports ITO PHTO 19.721399307250977 -155.04800415039062 Pacific/Honolulu airport OurAirports 3415 Hilo International Airport Hilo United States 38 -10 3416 ORL KORL 113 Orlando Executive Airport Orlando United States 28.5455 -81.332901 America/New_York airport OurAirports 66.91390228 647 A America/Anchorage airport OurAirports 6 3417 **Bettles Airport** Bettles United States BTT PABT -151.529007 -9 Clear Airport 552 3418 Clear Mews United States \N PACL 64.301201 -149.119995 America/Anchorage airport OurAirports Indian Mountain LRRS Airport Indian Mountains United States UTO PAIM 1273 3419 65.99279785 -153.7039948 A America/Anchorage airport OurAirports 3420 Fort Yukon Airport Fort Yukon United States FYU PFYU 66.57150268554688 -145.25 433 A America/Anchorage airport OurAirports Find the list of airlines having X stops

AirlineID SourceAirport SourceAirportID DestinationAirport DestinationAirportID Codeshare Stops Equipment Out[4]: Airline 5T 1623 YRT 132 YEK 50 ATR None 1 AC 330 ABJ 253 BRU 302 333 None YVR YBL AC 330 156 30 1 BEH None CU 1936 FCO 1555 HAV 1909 767 None FL HOU 3566 SAT 1316 3621 None 1 735 1316 MCO 3878 HOU 3566 None 73W ORF FL 1316 MCO 3878 717 3611 1 None ARN 737 **GEV** 715 ATP None 3448 MCO WN 4547 BOS 3878 73W None 1 WN 4547 MCO 3878 BOS 3448 None 73W Find list of airlines operating with code share (only showing top 10) There are two categories for "Codeshare" listed as NULL or "Y" routes_Codeshare = spark.sql(""" SELECT FROM routes WHERE Codeshare = 'Y'"") routes_Codeshare.toPandas().head(10).style Airline AirlinelD SourceAirport SourceAirportID DestinationAirport DestinationAirportID Codeshare Stops Equipment Out[5]: 0 2P 897 GES 2402 MNL 2397 0 320 2P 897 MNL 2397 GES 2402 0 320 4M 3201 DFW 3670 EZE 3988 Υ 0 777 DFW 3201 EZE 3670 0 777 4M 3988 EZE JFK 4M 3201 3988 3797 0 777 JFK 3797 EZE 777 4M 3201 3988 0 ARH **CSH** 5N 503 4362 6110 0 AN4 ARH 4362 MMK 5N 503 2949 AN4 USK 5N 503 ARH 4362 4369 0 AN4 CSH ARH 5N 503 6110 4362 AN4

Find the list of active airlines in the United Sates - Airline aggregation (only showing top 10) airlines_Active = spark.sql(""" FROM airlines WHERE Country = 'United States' AND Active = 'Y'"")

airlines_Active.toPandas().head(10).style AirlineID Name Alias ITAT ICAO Callsign **Country Active** Out[6]: 0 10 40-Mile Air Q5 MLA MILE-AIR United States 22 Aloha Airlines AQ AAH ALOHA United States 24 **American Airlines** AAL AMERICAN United States AA 35 G4 AAY ALLEGIANT United States 109 Alaska Central Express KO AER ACE AIR United States 210 Airlift International AIRLIFT United States 281 CACTUS United States America West Airlines HP AWE AIR WISCONSIN United States 282 Air Wisconsin ZW AWI 287 Allegheny Commuter Airlines ALLEGHENY United States \N None ALO

Find the country (or) territory has the highest number of Airports e.g. The table below shows the list of countries have the top 10 highest number of Airports airportsByCountry = spark.sql(""" SELECT Country, COUNT(DISTINCT AirportID) AS Airport_cnt FROM airports GROUP BY 1 ORDER BY 2 DESC""") airportsByCountry_top10 = airportsByCountry.toPandas().head(10) airportsByCountry_top10.style Out[7]: Country Airport_cnt 0 **United States** 1512

430 Canada 2 Australia 334 Russia 264 Brazil 264 Germany 249 China 241 217 France 8 United Kingdom 167 India 148 The horizontal barplot below can be used to visualize the "top 10 countries with the highest number of airports". It is easy to tell from the chart using either the width of the bars or the color of the bars. The country with the most airports, which is the United States, has the longest bar in the darkest color.

plt.style.use("ggplot") plt.figure(figsize = (12,8)) pal = sns.color_palette("Blues_d", len(airportsByCountry_top10["Airport_cnt"])) hbars = plt.barh(width="Airport_cnt", y="Country", data=airportsByCountry_top10, color=np.array(pal[::-1])) plt.gca().invert_yaxis() plt.title("Top 10 countries with the highest number of airports", fontsize = 20, weight = 'bold') plt.xlabel("# of Airports", size = 14) plt.ylabel("Country", size = 14) plt.bar_label(hbars, size = 12, weight = "bold") plt.show() Top 10 countries with the highest number of airports 1512 United States 430 Canada 334 Australia 249 Germany China 241 217 France 167 United Kingdom 148 India 1200

126

120

120

Coordinate Distance

Coordinate Distance

POINT (-118.334999 33.922798) 0.151406

POINT (-118.4079971 33.94250107) 0.190475

POINT (-118.035004 34.086102) 0.217492

POINT (-118.1520004 33.81769943) 0.249381

POINT (-118.450996399 34.0158004761) 0.204428

POINT (-118.33999633789 33.803398132324) 0.260435

POINT (-118.48999786377 34.209800720215) 0.290568

POINT (-118.413002014 34.2593002319) 0.267910

POINT (-117.980003357 33.8720016479) 0.321170

- 160

- 140

120

Airpor 80 5

60

- 40

- 20

POINT (-101.3740005 37.60400009)

POINT (-100.723999023 37.9275016785)

-100.88500213623047 POINT (-100.8850021362305 38.47430038452148)

-100.35600280761719 POINT (-100.3560028076172 37.27690124511719) 1.105112

-99.9655990600586 POINT (-99.96559906005859 37.76340103149414) 1.367376

34.20069885253906 -118.35900115966797 POINT (-118.359001159668 34.20069885253906) 0.188640

POINT (-100.8300018 37.49140167) 0.586726

POINT (-100.9599991 37.0442009) 0.837357

POINT (-101.508003235 36.6851005554) 1.122756

POINT (-102.68800354 38.0696983337) 1.383113

POINT (-101.879997 37.000702) 0.963873

111

104

100

100

92

88

83

ST_Point(CAST(Longitude as Double), CAST(Latitude as Double)) as Point

ST_Distance(Point, ST_Point(-101.332631, 37.794076)) as Distance

Longitude

-101.3740005

-100.8300018

-100.9599991

-101.879997

-101.508003235

-102.68800354

-100.723999023

80

of Incoming Airlines

Example 1: Find the top 10 closest airport to the point with the coordinate (37.794076, -101.332631)

Find the closest airport to a city X's geospatial coordinate

FROM airports""")

City, Latitude, Longitude,

Point as Coordinate,

Latitude

37.60400009

37.49140167

37.0442009

37.000702

Example 2: Find the top 10 closest airport to the point with the coordinate (34.04743, -118.24903)

ST_Distance(Point, ST_Point(-118.24903,34.04743)) as Distance

Latitude

33.922798

33.94250107

34.086102

33.81769943

33.803398132324

34.209800720215

WHERE ST_Contains(states.s_bound, airports_point.Point)

COUNT(DISTINCT AirportID) as airport_cnt

180

114

100

57

39

legend_kwds = {"label": "# of Airports", "shrink": 0.4})

fontsize = 14) if x["airport_cnt"] > 50 else False, axis = 1)

Georgia

Florida

34.2593002319

33.8720016479

Longitude

-118.334999

-118.4079971

-118.450996399

-118.035004

-118.1520004

-118.33999633789

-118.48999786377

ST_Point(CAST(Longitude as Double), CAST(Latitude as Double)) as Point

-118.413002014

-117.980003357

37.9275016785

38.474300384521484

37.27690124511719

Lamar 38.069698333699996

City

Hawthorne

Los Angeles

Santa Monica Municipal Airport Santa Monica 34.015800476100004

Find the airport in each US state's geospatial boundary

FROM airports""")

SELECT states.s_name, states.s_bound,

ORDER BY states.s_name""")

We can count the number of airports in each state and visualize using choropleth maps.

airports_point.AirportID, airports_point.Name, airports_point.City, airports_point.Point FROM states, airports_point

SELECT s_name as State,

FROM airport_per_state

21/11/25 01:28:19 WARN JoinQuery: UseIndex is true, but no index exists. Will build index on the fly.

airport_cnt_per_state_gpd.plot("airport_cnt", cmap = 'OrRd', figsize = (18,15), edgecolor = "0.8",

xy = x.s_bound.centroid.coords[0],

of airports in each US state

plt.title("# of airports in each US state", fontsize = 20, weight = 'bold')

plt.annotate(

text = x["State"],

s_bound airport_cnt

ORDER BY 3 DESC""")

airport_cnt_per_state_gpd = gpd.GeoDataFrame(airport_cnt_per_state_pd, geometry='s_bound')

s_bound,

GROUP BY 1,2

SELECT *,

airport_per_state.createOrReplaceTempView("airport_per_state")

The state having the most number of airports has the darkest color.

airport_cnt_per_state_pd = airport_cnt_per_state.toPandas()

Alaska POLYGON ((-141.02050 70.01870, -141.72910 70.1...

Texas POLYGON ((-106.57150 31.86590, -106.50420 31.7...

Florida POLYGON ((-87.60500 30.99880, -86.56130 30.996... Georgia POLYGON ((-85.60820 34.99740, -84.72660 34.990...

Michigan POLYGON ((-88.37130 48.30330, -87.60500 48.010...

Kansas POLYGON ((-102.05060 40.00340, -102.05060 40.0...

New York POLYGON ((-79.76240 42.51420, -79.06720 42.778... Wisconsin POLYGON ((-90.63030 42.51160, -87.02410 42.495...

6 North Carolina POLYGON ((-78.48500 33.79630, -79.67420 34.803...

airport_cnt_per_state_gpd.apply(lambda x:

California POLYGON ((-124.40090 41.99830, -123.62370 42.0...

airports_point.createOrReplaceTempView("airports_point")

Burbank

El Monte

Torrance

Van Nuys

Fullerton

Long Beach

Los Angeles

36.6851005554

FROM airports_point ORDER BY Distance""")

City

Ulysses

Sublette

Scott City

Liberal

Elkhart

Meade

Guymon

City, Latitude, Longitude,

Point as Coordinate,

FROM airports_point ORDER BY Distance""")

Name

Bob Hope Airport

Zamperini Field

Whiteman Airport

Van Nuys Airport

Los Angeles International Airport

Long Beach /Daugherty Field/ Airport

San Gabriel Valley Airport

Fullerton Municipal Airport

Garden City

airports_point.createOrReplaceTempView("airports_point")

dist_to_Pullman_pd = dist_to_Pullman.toPandas()

Name

Ulysses Airport

of Airports

COUNT(DISTINCT r.AirlineID) AS IncomeAirline_cnt

ON r.DestinationAirportID = a.AirportID

Find the top K cities with most incoming airlines

Airlines with the same airlineIDs coming from different source airports are considered as one airline

FROM routes r

GROUP BY 1

IncomeAirlinesByCity_top10 = IncomeAirlinesByCity.toPandas().head(10)

LEFT JOIN airports a

ORDER BY 2 DESC""")

The horizontal barplot below can be used to visualize the "top 10 cities with most incoming airlines".

The city with the most incoming airlines has the bar with the greatest length or the one with the darkest color.

pal = sns.color_palette("Blues_d", len(IncomeAirlinesByCity_top10["IncomeAirline_cnt"]))

plt.title("Top 10 cities with the most incoming airlines", fontsize = 20, weight = 'bold')

hbars = plt.barh(width="IncomeAirline_cnt", y="City", data=IncomeAirlinesByCity_top10, color=np.array(pal[::-1]))

Top 10 cities with the most incoming airlines

It is easy to tell from the chart using either the width of the bars or the color of the bars.

HAVING a.City IS NOT NULL

e.g. Find the top 10 cities with most incoming airlines

IncomeAirlinesByCity = spark.sql("""

IncomeAirlinesByCity_top10.style

City IncomeAirline_cnt

126

120

111

104 100

92

88

84 84

83

plt.xlabel("# of Incoming Airlines", size = 14)

plt.bar_label(hbars, size = 12, weight = "bold")

20

airports_point = spark.sql("""

dist_to_Pullman = spark.sql("""

dist_to_Pullman_pd.head(10).style

Sublette Municipal Airport

Garden City Regional Airport

Scott City Municipal Airport

Elkhart Morton County Airport

Meade Municipal Airport

Guymon Municipal Airport

Lamar Municipal Airport

dist_to_LA.toPandas().head(10).style

0 Jack Northrop Field Hawthorne Municipal Airport

airports_point = spark.sql("""

airport_per_state = spark.sql("""

airport_cnt_per_state = spark.sql("""

airport_cnt_per_state_gpd.head(10)

Liberal Mid-America Regional Airport

dist_to_LA = spark.sql("""

Paris

London

Moscow

Bangkok

Frankfurt

Rome

Dubai Istanbul

plt.figure(figsize = (12,8))

plt.ylabel("City", size = 14)

plt.gca().invert_yaxis()

8 Hong Kong

9 Singapore

plt.show()

Paris

London

Moscow

Bangkok

Frankfurt

Rome

Dubai

Istanbul

Hong Kong

Singapore

In [11]:

In [12]:

Out[12]:

In [13]:

Out[13]:

In [14]:

In [15]:

In [16]:

Out[16]:

In [17]:

plt.axis("off")

plt.show()

0

Out[9]:

In [10]: