

```

1 # Run this cell to import pyspark and to define start_spark() and stop_spark()
2
3 import findspark
4
5 findspark.init()
6
7 import getpass
8 import random
9 import re
10
11 import pandas
12 import pyspark
13 from IPython.display import HTML, display
14 from pyspark import SparkContext
15 from pyspark.sql import SparkSession
16
17 # Constants used to interact with Azure Blob Storage using the hdfs command or Spark
18
19 global username
20
21 username = re.sub("@.*", "", getpass.getuser())
22
23 global azure_account_name
24 global azure_data_container_name
25 global azure_user_container_name
26 global azure_user_token
27
28 azure_account_name = "madsstorage002"
29 azure_data_container_name = "campus-data"
30 azure_user_container_name = "campus-user"
31 azure_user_token = r"sp=racwdl&st=2025-08-01T09:41:33Z&se=2026-12-30T16:56:33Z&spr=https&sv=2024-11-04&sr=c&sig=Gz"
32
33
34 # Functions used below
35
36
37 def dict_to_html(d):
38     """Convert a Python dictionary into a two column table for display."""
39
40     html = []
41
42     html.append(f'<table width="100%" style="width:100%; font-family: monospace;">')
43     for k, v in d.items():
44         html.append(f'<tr><td style="text-align:left;">{k}</td><td>{v}</td></tr>')
45     html.append(f'</table>')
46
47     return "".join(html)
48
49
50 def show_as_html(df, n=20):
51     """Leverage existing pandas jupyter integration to show a spark dataframe as html.
52
53     Args:
54         n (int): number of rows to show (default: 20)
55     """
56
57     display(df.limit(n).toPandas())
58
59
60 def display_spark():
61     """Display the status of the active Spark session if one is currently running."""
62
63     if "spark" in globals() and "sc" in globals():
64
65         name = sc.getConf().get("spark.app.name")
66
67         html = [
68             f"<p><b>Spark</b></p>",
69             f"<p>The spark session is <b><span style='color:green'>active</span></b>, look for <code>{name}</code></p>",
70             f"<ul>",
71             f"<li><a href='http://localhost:{sc.uiWebUrl.split(':')[1]}' target='_blank'>Spark Application UI</a></li></ul>",
72             f"<p><b>Config</b></p>",
73             dict_to_html(dict(sc.getConf().getAll())),
74             f"<p><b>Notes</b></p>",
75             f"<ul>",
76             f"<li>The spark session <code>spark</code> and spark context <code>sc</code> global variables have been created.",
77             f"<li>Please run <code>stop_spark()</code> before closing the notebook or restarting the kernel or killing the notebook.",
78             f"</ul>",
79         ]
80
81     display(HTML("".join(html)))

```

```

82
83     else:
84
85         html = [
86             f"<p><b>Spark</b></p>",
87             f"<p>The spark session is <b><span style='color:red'>stopped</span></b>, confirm that <code>{username}</code>",
88             f"<ul>",
89             f"<li><a href='http://mathmadslinux2p.canterbury.ac.nz:8080/' target='_blank'>Spark UI</a></li>",
90             f"</ul>",
91         ]
92         display(HTML("".join(html)))
93
94
95 # Functions to start and stop spark
96
97
98 def start_spark(
99     executor_instances=2, executor_cores=1, worker_memory=1, master_memory=1
100 ):
101     """Start a new Spark session and define globals for SparkSession (spark) and SparkContext (sc).
102
103     Args:
104         executor_instances (int): number of executors (default: 2)
105         executor_cores (int): number of cores per executor (default: 1)
106         worker_memory (float): worker memory (default: 1)
107         master_memory (float): master memory (default: 1)
108     """
109
110     global spark
111     global sc
112
113     cores = executor_instances * executor_cores
114     partitions = cores * 4
115     port = 4000 + random.randint(1, 999)
116
117     spark = (
118         SparkSession.builder.config(
119             "spark.driver.extraJavaOptions",
120             f"-Dderby.system.home=/tmp/{username}/spark/",
121         )
122         .config("spark.dynamicAllocation.enabled", "false")
123         .config("spark.executor.instances", str(executor_instances))
124         .config("spark.executor.cores", str(executor_cores))
125         .config("spark.cores.max", str(cores))
126         .config("spark.driver.memory", f"{master_memory}g")
127         .config("spark.executor.memory", f"{worker_memory}g")
128         .config("spark.driver.maxResultSize", "0")
129         .config("spark.sql.shuffle.partitions", str(partitions))
130         .config(
131             "spark.kubernetes.container.image",
132             "madsregistry001.azurecr.io/hadoop-spark:v3.3.5-openjdk-8",
133         )
134         .config("spark.kubernetes.container.image.pullPolicy", "IfNotPresent")
135         .config("spark.kubernetes.memoryOverheadFactor", "0.3")
136         .config("spark.memory.fraction", "0.1")
137         .config(
138             f"fs.azure.sas.{azure_user_container_name}.{azure_account_name}.blob.core.windows.net",
139             azure_user_token,
140         )
141         .config("spark.app.name", f"{username} (notebook)")
142         .getOrCreate()
143     )
144     sc = SparkContext.getOrCreate()
145
146     display_spark()
147
148
149 def stop_spark():
150     """Stop the active Spark session and delete globals for SparkSession (spark) and SparkContext (sc)."""
151
152     global spark
153     global sc
154
155     if "spark" in globals() and "sc" in globals():
156
157         spark.stop()
158
159         del spark
160         del sc
161
162     display_spark()
163
164

```

```

165 # Make css changes to improve spark output readability
166
167 html = [
168     "<style>",
169     "pre { white-space: pre !important; }",
170     "table.dataframe td { white-space: nowrap !important; }",
171     "table.dataframe thead th:first-child, table.dataframe tbody th { display: none; }",
172     "</style>",
173 ]
174 display(HTML("".join(html)))

```

```

1 # Run this cell to start a spark session in this notebook
2
3 start_spark(executor_instances=4, executor_cores=2, worker_memory=4, master_memory=4)

```

```

1 Warning: Ignoring non-Spark config property: fs.azure.sas.campus-user.madsstorage002.blob.core.windows.net
2 Warning: Ignoring non-Spark config property: SPARK_DRIVER_BIND_ADDRESS
3 25/09/11 22:13:42 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java
4 Setting default log level to "WARN".
5 To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).

```

## Spark

The spark session is **active**, look for `yxi75 (notebook)` under the running applications section in the Spark UI.

- [Spark Application UI](#)

## Config

spark.dynamicAllocation.enabled	false
spark.fs.azure.sas.uco-user.madsstorage002.blob.core.windows.net	"sp=racwdl&st=2024-09-19T08:00:18Z&se=2025-09-19T16:00:18Z&spr=https&sv=2022-11-02&sr=c&sig=qtg6fCdoFz6k3EJLw7dA8D3D8wN0neAYw8yG4z4Lw2o%3D"
spark.kubernetes.driver.pod.name	spark-master-driver
spark.executor.instances	4
spark.app.id	spark-1054e131d02d4bc6a6e42e1ed4b1e027
spark.kubernetes.executor.podNamePrefix	yxi75-notebook-344ce89938444bdb
spark.driver.memory	4g
spark.app.name	yxi75 (notebook)
spark.fs.azure.sas.campus-user.madsstorage002.blob.core.windows.net	"sp=racwdl&st=2024-09-19T08:03:31Z&se=2025-09-19T16:03:31Z&spr=https&sv=2022-11-02&sr=c&sig=kMP%2BsBsRzdVVR8rrg%2BNbDhkRBNS6Q98kYY695XMRFDU%3D"
spark.kubernetes.container.image.pullPolicy	IfNotPresent
spark.sql.shuffle.partitions	32
spark.kubernetes.namespace	yxi75
spark.serializer.objectStreamReset	100
spark.driver.maxResultSize	0
spark.app.submitTime	1757585623019
spark.submit.deployMode	client
spark.master	k8s://https://kubernetes.default.svc.cluster.local:443
spark.driver.extraJavaOptions	-Djava.net.preferIPv6Addresses=false --add-opens=java.base/java.lang=ALL-UNNAMED --add-opens=java.base/java.lang.invoke=ALL-UNNAMED --add-opens=java.base/java.lang.reflect=ALL-UNNAMED --add-opens=java.base/java.io=ALL-UNNAMED --add-opens=java.base/java.net=ALL-UNNAMED --add-opens=java.base/java.nio=ALL-UNNAMED --add-opens=java.base/java.util=ALL-UNNAMED --add-opens=java.base/java.util.concurrent=ALL-UNNAMED --add-opens=java.base/java.util.concurrent.atomic=ALL-UNNAMED --add-opens=java.base/jdk.internal.ref=ALL-UNNAMED --add-

	opens=java.base/sun.nio.ch=ALL-UNNAMED --add- opens=java.base/sun.nio.cs=ALL-UNNAMED --add- opens=java.base/sun.security.action=ALL-UNNAMED --add- opens=java.base/sun.util.calendar=ALL-UNNAMED --add- opens=java.security.jgss/sun.security.krb5=ALL-UNNAMED - Djdk.reflect.useDirectMethodHandle=false - Dderby.system.home=/tmp/yxi75/spark/
spark.fs.azure	org.apache.hadoop.fs.azure.NativeAzureFileSystem
spark.app.startTime	1757585623186
spark.memory.fraction	0.1
spark.executor.memory	4g
spark.executor.id	driver
spark.kubernetes.executor.container.image	madsregistry001.azurecr.io/hadoop-spark:v3.3.5-openjdk-8-1.0.16
spark.executor.cores	2
spark.kubernetes.memoryOverheadFactor	0.3
spark.driver.host	spark-master-svc
spark.ui.port	\${env:SPARK_UI_PORT}
spark.kubernetes.container.image	madsregistry001.azurecr.io/hadoop-spark:v3.3.5-openjdk-8
spark.kubernetes.executor.podTemplateFile	/opt/spark/conf/executor-pod-template.yaml
fs.azure.sas.campus-user.madsstorage002.blob.core.windows.net	sp=racwdl&st=2025-08-01T09:41:33Z&se=2026-12- 30T16:56:33Z&spr=https&sv=2024-11- 04&sr=c&sig=GzR1hq7EJ01RHj92oD01MBNjkc602nrpfB5H8C17FFY%3D
spark.rdd.compress	True
spark.executor.extraJavaOptions	-Djava.net.preferIPv6Addresses=false - XX:+IgnoreUnrecognizedVMOptions --add- opens=java.base/java.lang=ALL-UNNAMED --add- opens=java.base/java.lang.invoke=ALL-UNNAMED --add- opens=java.base/java.lang.reflect=ALL-UNNAMED --add- opens=java.base/java.io=ALL-UNNAMED --add- opens=java.base/java.net=ALL-UNNAMED --add- opens=java.base/java.nio=ALL-UNNAMED --add- opens=java.base/java.util=ALL-UNNAMED --add- opens=java.base/java.util.concurrent=ALL-UNNAMED --add- opens=java.base/java.util.concurrent.atomic=ALL-UNNAMED --add- opens=java.base/jdk.internal.ref=ALL-UNNAMED --add- opens=java.base/sun.nio.ch=ALL-UNNAMED --add- opens=java.base/sun.nio.cs=ALL-UNNAMED --add- opens=java.base/sun.security.action=ALL-UNNAMED --add- opens=java.base/sun.util.calendar=ALL-UNNAMED --add- opens=java.security.jgss/sun.security.krb5=ALL-UNNAMED - Djdk.reflect.useDirectMethodHandle=false
spark.cores.max	8
spark.driver.port	7077
spark.submit.pyFiles	
spark.ui.showConsoleProgress	true

## Notes

- The spark session `spark` and spark context `sc` global variables have been defined by `start_spark()`.
- Please run `stop_spark()` before closing the notebook or restarting the kernel or kill `yxi75 (notebook)` by hand using the link in the Spark UI.

```

1  # Write your imports here or insert cells below
2
3  import re
4  import subprocess
5  from math import asin, cos, radians, sin, sqrt
6  from pprint import pprint
7
8  import numpy as np
9  import pandas as pd
10 from pyspark.sql import functions as F
11 from pyspark.sql.functions import udf
12 from pyspark.sql.types import *

```

```

1  # Paths global variables
2  DATA_ROOT = "wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd/"
3  USER_ROOT = "wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/"
4
5  paths = {
6      "daily": DATA_ROOT + "daily/",
7      "stations": DATA_ROOT + "ghcnd-stations.txt",
8      "countries": DATA_ROOT + "ghcnd-countries.txt",
9      "states": DATA_ROOT + "ghcnd-states.txt",
10     "inventory": DATA_ROOT + "ghcnd-inventory.txt",
11 }
12
13 paths
14
15 stations_enriched_savepath = USER_ROOT + "stations_enriched_parquet"
16 station_count_by_country_path = USER_ROOT + "station_count_by_country_parquet"
17 station_count_us_terri_path = USER_ROOT + "station_count_us_terri_parquet"
18 country_meta_with_station_num_path = USER_ROOT + "country_meta_with_station_num"
19 states_meta_with_station_num_path = USER_ROOT + "states_meta_with_station_num"

```

## Analysis

### Q1 Station Study

(a)

```

1  # load stations_enriched from saved path
2  stations_enriched = spark.read.parquet(stations_enriched_savepath)
3  stations_enriched.cache()
4  # check variable
5  show_as_html(stations_enriched)

```

```

1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }

```

ID	STATE	COUNTRY_CODE	LATITUDE	LONGITUDE	ELEVATION	NAME	GSN_FLAG	HCN
AE000041196		AE	25.3330	55.5170	34.0	SHARJAH INTER. AIRP	GSN	
AEM00041218		AE	24.2620	55.6090	264.9	AL AIN INTL		
AGE00147715		AG	35.4200	8.1197	863.0	TEBESSA		
AGE00147717		AG	35.2000	0.6300	476.0	SIDI-BEL-ABBES		
AGE00147719		AG	33.7997	2.8900	767.0	LAGHOUAT		
AGM00060425		AG	36.2130	1.3320	141.1	ECH CHELIFF		
AGM00060430		AG	36.3000	2.2330	721.0	MILIANA		
AGM00060490		AG	35.6240	-0.6210	89.9	ES SENIA		
AGM00060514		AG	35.1670	2.3170	801.0	KSAR CHELLALA		
AGM00060515		AG	35.3330	4.2060	459.0	BOU SAADA		
AGM00060522		AG	34.8200	-1.7700	426.0	MAGHNIA		
AGM00060536		AG	34.8670	0.1500	752.0	SAIDA		
AGM00060555		AG	33.0680	6.0890	85.0	SIDI MAHDI		
AGM00060580		AG	31.9170	5.4130	150.0	OUARGLA		
AJ000037742		AJ	40.9000	47.3000	313.0	ORDJONIKIDZE,ZERNOSOVHOZ		
AJ000037747		AJ	40.6170	47.1500	15.0	EVLAKH AIRPORT		
AJ000037756		AJ	40.5330	48.9330	755.0	MARAZA		
AJ000037816		AJ	40.5000	46.1000	1655.0	DASHKESAN		
AJ000037899		AJ	39.7670	46.7500	1355.0	SHUSHA		
AJ000037914		AJ	39.9000	48.0000	-1.0	IMISLY		

20 rows × 21 columns

```
1 # How many stations are there in total?
2 total_stations = stations_enriched.count()
3 print("total_stations: ", total_stations)
4
5 # How many stations were active so far in 2025
6 active_2025 = stations_enriched.filter(F.col("LASTYEAR_ANY") >= 2025).count()
7 print("active stations in 2025: ", active_2025)
```

```
1 total_stations: 129657
2 active stations in 2025: 38481
```

```
1 # station network belongings analysis
2 network_counts = stations_enriched.select(
3     F.when(F.col("GSN_FLAG").contains("GSN"), 1).otherwise(0).alias("is_GSN"),
4     F.when(F.col("HCN_CRN").contains("HCN"), 1).otherwise(0).alias("is_HCN"),
5     F.when(F.col("HCN_CRN").contains("CRN"), 1).otherwise(0).alias("is_CRN"),
6 ).agg(
7     F.sum("is_GSN").alias("GSN"),
8     F.sum("is_HCN").alias("HCN"),
9     F.sum("is_CRN").alias("CRN"),
10 )
11
12 show_as_html(network_counts)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	GSN	HCN	CRN
991	1218	234	

```

1 # Are there any stations that are in more than one of these networks?
2 # multi-network station counts query criteria: net_count >= 2
3 flags = stations_enriched.select(
4     "ID",
5     F.when(F.col("GSN_FLAG").contains("GSN"), 1).otherwise(0).alias("is_GSN"),
6     F.when(F.col("HCN_CRN").contains("HCN"), 1).otherwise(0).alias("is_HCN"),
7     F.when(F.col("HCN_CRN").contains("CRN"), 1).otherwise(0).alias("is_CRN"),
8 )
9
10 multi_network = (
11     flags.withColumn("net_count", F.col("is_GSN") + F.col("is_HCN") + F.col("is_CRN"))
12     .filter(F.col("net_count") >= 2)
13     .select("ID", "is_GSN", "is_HCN", "is_CRN", "net_count")
14 )
15 show_as_html(multi_network)
16
17 print(
18     f"multi_network number is {multi_network.count()}, and they all overlap with GSN and HCN."
19 )

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe tthead th {
6     text-align: right;
7 }

```

ID	is_GSN	is_HCN	is_CRN	net_count
USW00012921	1	1	0	2
USW00024128	1	1	0	2
USW00024213	1	1	0	2
USW00012836	1	1	0	2
USW00093193	1	1	0	2
USW00014771	1	1	0	2
USW00014922	1	1	0	2
USW00094008	1	1	0	2
USW00024144	1	1	0	2
USW00003870	1	1	0	2
USW00023044	1	1	0	2
USW00013782	1	1	0	2
USW00023051	1	1	0	2
USW00014742	1	1	0	2
USW00093729	1	1	0	2

```

1 | multi_network number is 15, and they all overlap with GSN and HCN.

```

(b)

```

1 # How many stations are there in the Southern Hemisphere?
2 southern_hemisphere = stations_enriched.filter(F.col("LATITUDE") < 0).count()
3 print("southern_hemisphere count: ", southern_hemisphere)

```

```

1 | southern_hemisphere count: 25357

```

```

1 # identify us_territories, excluding the United States it self
2 # How many stations are there in total in the territories of the United States, excluding the United States itself?
3 us_territories = (
4     stations_enriched.filter(
5         (F.col("COUNTRY_NAME").contains("United States"))
6         & (F.trim(F.col("COUNTRY_NAME")) != "United States")

```

```
7     )
8     .groupBy("COUNTRY_NAME")
9     .agg(F.count("*").alias("STATION_NUM"))
10  )
11
12
13  show_as_html(us_territories)
14  us_territories_station_num = us_territories.agg(
15      F.sum("station_num").alias("us_territories_station_num")
16  )
17  show_as_html(us_territories_station_num)
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

COUNTRY_NAME	STATION_NUM
Northern Mariana Islands [United States]	11
Puerto Rico [United States]	260
Guam [United States]	34
Johnston Atoll [United States]	4
American Samoa [United States]	21
Virgin Islands [United States]	77
Midway Islands [United States]	3
Palmyra Atoll [United States]	3
Wake Island [United States]	1

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

us_territories_station_num
414

(c)

```
1  # groupby country and states and count, save
2  station_num_by_country = (
3      stations_enriched.groupBy("COUNTRY_CODE", "COUNTRY_NAME")
4      .agg(F.count("*").alias("STATION_NUM"))
5      .orderBy(F.desc("STATION_NUM"))
6  )
7
8  show_as_html(station_num_by_country)
```

```
1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }
```

COUNTRY_CODE	COUNTRY_NAME	STATION_NUM
US	United States	75846
AS	Australia	17088
CA	Canada	9269



BR	Brazil	5989
MX	Mexico	5249
IN	India	3807
SW	Sweden	1721
SF	South Africa	1166
RS	Russia	1123
GM	Germany	1123
FI	Finland	922
NO	Norway	461
NL	Netherlands	386
KZ	Kazakhstan	329
WA	Namibia	283
RQ	Puerto Rico [United States]	260
CH	China	228
SP	Spain	207
UP	Ukraine	204
JA	Japan	202

```
1 us_territories.write.parquet(station_count_us_terri_path, mode="overwrite")
2 station_count_by_country.write.parquet(station_count_by_country_path, mode="overwrite")
```

```
1 # check save result
2 !hdfs dfs -ls -h {station_count_us_terri_path}
3 !hdfs dfs -ls -h {station_count_by_country_path}
```

```
1 Found 2 items
2 -rw-r--r-- 1 yxi75 supergroup 0 2025-09-11 22:33 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/station_count_us_te
3 -rw-r--r-- 1 yxi75 supergroup 1.0 K 2025-09-11 22:33 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/station_count_us_te
4 Found 2 items
5 -rw-r--r-- 1 yxi75 supergroup 0 2025-09-11 22:33 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/station_count_by_co
6 -rw-r--r-- 1 yxi75 supergroup 4.9 K 2025-09-11 22:33 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/station_count_by_co
```

```
1 # load countries meta table for join
2 countries_df = (
3     spark.read.text(paths["countries"])
4     .withColumn("CODE", F.substring("value", 1, 2))
5     .withColumn("COUNTRY_NAME", F.substring("value", 4, 61))
6     .withColumnRenamed("CODE", "COUNTRY_CODE")
7     .drop("value")
8 )
9
10 # join country counts onto countries
11 country_meta_with_station_num = countries_df.join(
12     station_count_by_country.withColumnRenamed("count", "STATION_NUM"),
13     on=["COUNTRY_NAME", "COUNTRY_CODE"],
14     how="left",
15 ).orderBy("COUNTRY_NAME")
16
17 show_as_html(country_meta_with_station_num)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

COUNTRY_NAME	COUNTRY_CODE	STATION_NUM
Afghanistan	AF	4

Albania	AL	3
Algeria	AG	87
American Samoa [United States]	AQ	21
Angola	AO	6
Antarctica	AY	102
Antigua and Barbuda	AC	2
Argentina	AR	101
Armenia	AM	53
Australia	AS	17088
Austria	AU	13
Azerbaijan	AJ	66
Bahamas, The	BF	72
Bahrain	BA	1
Bangladesh	BG	10
Barbados	BB	1
Belarus	BO	51
Belgium	BE	1
Belize	BH	1
Benin	BN	9

```
1 # save to country_meta_with_station_num_path
2 country_meta_with_station_num.write.parquet(
3     country_meta_with_station_num_path, "overwrite"
4 )
```

```
1 25/09/11 22:55:21 WARN AzureFileSystemThreadPoolExecutor: Disabling threads for Delete operation as thread count 0 is <= 1
```

```
1 !hdfs dfs -ls {country_meta_with_station_num_path}
```

```
1 Found 2 items
2 -rw-r--r--  1 yxi75 supergroup          0 2025-09-11 22:55 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/country_meta_with_s
3 -rw-r--r--  1 yxi75 supergroup    4945 2025-09-11 22:55 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/country_meta_with_s
```

```
1 station_count_by_state = stations_enriched.groupBy("STATE", "STATE_NAME").agg(
2     F.count("*").alias("STATION_NUM")
3 )
4 show_as_html(station_count_by_state, 3)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

STATE	STATE_NAME	STATION_NUM
NS	NOVA SCOTIA	398
NE	NEBRASKA	2436
IA	IOWA	1106

```
1 states_df = (  
2     spark.read.text(paths["states"])  
3     .withColumn("CODE", F.substring("value", 1, 2))  
4     .withColumn("STATE_NAME", F.substring("value", 4, 47))  
5     .drop("value")  
6 )  
7 show_as_html(states_df, 3)
```

```
1 .dataframe tbody tr th {  
2     vertical-align: top;  
3 }  
4  
5 .dataframe thead th {  
6     text-align: right;  
7 }
```

CODE	STATE_NAME
AB	ALBERTA
AK	ALASKA
AL	ALABAMA

```
1 states_meta_with_station_num = (  
2     states_df.withColumnRenamed("CODE", "STATE_CODE")  
3     .join(  
4         station_count_by_state.withColumnRenamed("STATE", "STATE_CODE"),  
5         on=["STATE_NAME", "STATE_CODE"],  
6         how="left",  
7     )  
8     .orderBy(F.desc("STATION_NUM"))  
9 )  
10  
11 show_as_html(states_meta_with_station_num)
```

```
1 .dataframe tbody tr th {  
2     vertical-align: top;  
3 }  
4  
5 .dataframe thead th {  
6     text-align: right;  
7 }
```

STATE_NAME	STATE_CODE	STATION_NUM
TEXAS	TX	6472
COLORADO	CO	4784
CALIFORNIA	CA	3166
NORTH CAROLINA	NC	2747
MINNESOTA	MN	2675
NEBRASKA	NE	2436
KANSAS	KS	2401
NEW MEXICO	NM	2295
FLORIDA	FL	2244
ILLINOIS	IL	2234
OREGON	OR	2031
ONTARIO	ON	2021
INDIANA	IN	2020
NEW YORK	NY	1912
TENNESSEE	TN	1755
BRITISH COLUMBIA	BC	1720
WASHINGTON	WA	1694

ARIZONA	AZ	1692
PENNSYLVANIA	PA	1641
MISSOURI	MO	1624

```
1 states_meta_with_station_num.write.parquet(  
2     states_meta_with_station_num_path, "overwrite"  
3 )
```

```
1 | 25/09/11 23:18:14 WARN AzureFileSystemThreadPoolExecutor: Disabling threads for Delete operation as thread count 0 is <= 1
```

```
1 # check USER_ROOT storage status  
2 !hdfs dfs -ls -h {USER_ROOT}
```

```
1 Found 5 items  
2 drwxr-xr-x - yxi75 supergroup      0 2025-09-11 22:55 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/country_meta_with_s  
3 drwxr-xr-x - yxi75 supergroup      0 2025-09-11 23:18 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/states_meta_with_st  
4 drwxr-xr-x - yxi75 supergroup      0 2025-09-11 22:33 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/station_count_by_co  
5 drwxr-xr-x - yxi75 supergroup      0 2025-09-11 22:33 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/station_count_us_te  
6 drwxr-xr-x - yxi75 supergroup      0 2025-09-11 18:30 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/stations_enriched_p
```

## Q2 Distance User Define Function

### (a) User Define Function: Haversine Distance

#### Haversine Distance Formula

The **Haversine formula** is used to calculate the great-circle distance between two points on a sphere (such as the Earth), given their latitude and longitude.

### 1. Definitions

Let two points be defined as:

- Point 1:  $(\phi_1, \lambda_1)$
- Point 2:  $(\phi_2, \lambda_2)$

where:

- $\phi$  = latitude (in radians)
- $\lambda$  = longitude (in radians)
- $R$  = radius of the Earth (mean radius  $\approx 6,371$  km)

### 2. Formula

The haversine function is defined as:

$$\text{hav}(\theta) = \sin^2\left(\frac{\theta}{2}\right) \tag{1}$$

Using this, the central angle  $c$  between two points is:

$$a = \sin^2\left(\frac{\Delta\phi}{2}\right) + \cos(\phi_1) \cdot \cos(\phi_2) \cdot \sin^2\left(\frac{\Delta\lambda}{2}\right) \tag{2}$$

$$c = 2 \cdot \arcsin(\sqrt{a}) \tag{3}$$

Finally, the distance  $d$  is:

$$d = R \cdot c \tag{4}$$

where:

- $\Delta\phi = \phi_2 - \phi_1$
- $\Delta\lambda = \lambda_2 - \lambda_1$

### 3. Notes

- The Haversine formula accounts for Earth's curvature, making it more accurate than a simple Euclidean distance in geographic applications.
- It is sufficiently precise for most data science and GIS tasks.

- For centimeter-level accuracy (e.g., geodesy, GPS), ellipsoidal models such as **Vincenty's formula** or the **WGS84 geodesic** method should be used.

(This part is generated by chatGPT)

Haversine UDF

```
1 # Haversine UDF
2 from math import asin, cos, radians, sin, sqrt
3
4 from pyspark.sql.functions import udf
5
6
7 def haversine_km(lat1, lon1, lat2, lon2):
8     R = 6371.0088 # mean Earth radius in km
9     phi1, phi2 = radians(lat1), radians(lat2)
10    dphi = radians(lat2 - lat1)
11    dlamba = radians(lon2 - lon1)
12    a = sin(dphi / 2) ** 2 + cos(phi1) * cos(phi2) * sin(dlamba / 2) ** 2
13    c = 2 * asin(sqrt(a))
14    return float(R * c)
15
16
17 # register UDF
18 haversine_udf = udf(haversine_km, DoubleType())
```

test udf within subset stations

```
1 # test udf within subset stations
2 test_stations = stations_enriched.limit(10)
3 show_as_html(test_stations)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

ID	STATE	COUNTRY_CODE	LATITUDE	LONGITUDE	ELEVATION	NAME	GSN_FLAG	HCN_CRN	WMO_ID	..
AE000041196		AE	25.3330	55.5170	34.0	SHARJAH INTER. AIRP	GSN		41196	..
AEM00041218		AE	24.2620	55.6090	264.9	AL AIN INTL			41218	..
AGE00147715		AG	35.4200	8.1197	863.0	TEBESSA				..
AGE00147717		AG	35.2000	0.6300	476.0	SIDI-BEL-ABBES				..
AGE00147719		AG	33.7997	2.8900	767.0	LAGHOUAT			60545	..
AGM00060425		AG	36.2130	1.3320	141.1	ECH CHELIFF			60425	..
AGM00060430		AG	36.3000	2.2330	721.0	MILIANA			60430	..
AGM00060490		AG	35.6240	-0.6210	89.9	ES SENIA			60490	..
AGM00060514		AG	35.1670	2.3170	801.0	KSAR CHELLALA			60514	..
AGM00060515		AG	35.3330	4.2060	459.0	BOU SAADA			60515	..

10 rows × 21 columns

```
1 # generate station pairs
2
3 # SELF CROSS JOIN within test_stations
4 left = test_stations.select(
5     F.col("ID").alias("ID_A"),
6     F.col("NAME").alias("NAME_A"),
7     F.col("LATITUDE").alias("LAT_A"),
8     F.col("LONGITUDE").alias("LON_A"),
9 )
10 right = test_stations.select(
11     F.col("ID").alias("ID_B"),
12     F.col("NAME").alias("NAME_B"),
13     F.col("LATITUDE").alias("LAT_B"),
14     F.col("LONGITUDE").alias("LON_B"),
15 )
```

```
16
17 test_pairs = left.crossJoin(right).filter(F.col("ID_A") < F.col("ID_B"))
18
19 test_pairs = test_pairs.withColumn(
20     "DIST_KM", haversine_udf("LAT_A", "LON_A", "LAT_B", "LON_B")
21 )
22 show_as_html(test_pairs)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

ID_A	NAME_A	LAT_A	LON_A	ID_B	NAME_B	LAT_B	LON_B	DIST_KM
AG000060680	TAMANRASSET	22.800	5.4331	AGE00147794	BEJAIA-CAP CARBON	36.780	5.100	1554.836252
AG000060680	TAMANRASSET	22.800	5.4331	AGM00060351	JIJEL	36.795	5.874	1556.750841
AG000060680	TAMANRASSET	22.800	5.4331	AGM00060353	JIJEL-PORT	36.817	5.883	1559.219777
AG000060680	TAMANRASSET	22.800	5.4331	AGM00060419	MOHAMED BOUDIAF INTL	36.276	6.620	1502.817981
AG000060680	TAMANRASSET	22.800	5.4331	AGM00060550	EL-BAYADH	33.667	1.000	1283.611711
AG000060680	TAMANRASSET	22.800	5.4331	AGM00060563	HASSIR'MEL	32.933	3.283	1146.297735
AG000060680	TAMANRASSET	22.800	5.4331	AGM00060670	TISKA	24.293	9.452	442.002990
AG000060680	TAMANRASSET	22.800	5.4331	AJ000037639	AGSTAPHA AIRPORT	41.133	45.417	4236.617857
AG000060680	TAMANRASSET	22.800	5.4331	AJ000037656	ADJINAURSKAYA_STEP'	41.200	46.800	4350.094399
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AGM00060351	JIJEL	36.795	5.874	68.946174
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AGM00060353	JIJEL-PORT	36.817	5.883	69.838733
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AGM00060419	MOHAMED BOUDIAF INTL	36.276	6.620	146.921796
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AGM00060550	EL-BAYADH	33.667	1.000	508.348378
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AGM00060563	HASSIR'MEL	32.933	3.283	458.743957
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AGM00060670	TISKA	24.293	9.452	1449.209308
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AJ000037639	AGSTAPHA AIRPORT	41.133	45.417	3488.136797
AGE00147794	BEJAIA-CAP CARBON	36.780	5.1000	AJ000037656	ADJINAURSKAYA_STEP'	41.200	46.800	3602.576212
AGM00060351	JIJEL	36.795	5.8740	AGM00060353	JIJEL-PORT	36.817	5.883	2.574176
AGM00060351	JIJEL	36.795	5.8740	AGM00060419	MOHAMED BOUDIAF INTL	36.276	6.620	88.162740
AGM00060351	JIJEL	36.795	5.8740	AGM00060550	EL-BAYADH	33.667	1.000	562.845643

```
1 # select New Zealand stations
2 nz_stations = stations_enriched.filter(
3     F.col("COUNTRY_NAME").contains("New Zealand")
4 ).select("ID", "NAME", "LATITUDE", "LONGITUDE")
5
6 # SELF CROSS JOIN within newzealand
7 left = nz_stations.select(
8     F.col("ID").alias("ID_A"),
9     F.col("NAME").alias("NAME_A"),
10    F.col("LATITUDE").alias("LAT_A"),
11    F.col("LONGITUDE").alias("LON_A"),
12 )
13 right = nz_stations.select(
14     F.col("ID").alias("ID_B"),
15     F.col("NAME").alias("NAME_B"),
16     F.col("LATITUDE").alias("LAT_B"),
17     F.col("LONGITUDE").alias("LON_B"),
18 )
19
20 pairs = left.crossJoin(right).filter(F.col("ID_A") < F.col("ID_B"))
21
22 pairs = pairs.withColumn("DIST_KM", haversine_udf("LAT_A", "LON_A", "LAT_B", "LON_B"))
```

```
1 | show_as_html(pairs)
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

ID_A	NAME_A	LAT_A	LON_A	ID_B	NAME_B	LAT_B	LON_B	DIST_KM
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZ000937470	TARA HILLS	-44.517	169.900	218.309321
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZ000939870	CHATHAM ISLANDS AWS	-43.950	-176.567	1015.251566
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZM00093781	CHRISTCHURCH INTL	-43.489	172.532	152.258567
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZM00093110	AUCKLAND AERO AWS	-37.000	174.800	714.127832
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	TL000091724	NUKUNONO	-9.200	-171.917	4082.088199
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZ000939450	CAMPBELL ISLAND AWS	-52.550	169.167	1101.720586
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZM00093439	WELLINGTON AERO AWS	-41.333	174.800	350.796507
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZM00093929	ENDERBY ISLAND AWS	-50.483	166.300	934.248833
NZ000936150	HOKITIKA AERODROME	-42.717	170.983	NZM00093678	KAIKOURA	-42.417	173.700	224.981589
NZ000937470	TARA HILLS	-44.517	169.900	NZ000939870	CHATHAM ISLANDS AWS	-43.950	-176.567	1078.799953
NZ000937470	TARA HILLS	-44.517	169.900	NZM00093781	CHRISTCHURCH INTL	-43.489	172.532	239.530461
NZ000939870	CHATHAM ISLANDS AWS	-43.950	-176.567	NZM00093781	CHRISTCHURCH INTL	-43.489	172.532	876.909075
NZ000937470	TARA HILLS	-44.517	169.900	NZM00093110	AUCKLAND AERO AWS	-37.000	174.800	931.744249
NZ000939870	CHATHAM ISLANDS AWS	-43.950	-176.567	NZM00093110	AUCKLAND AERO AWS	-37.000	174.800	1062.046480
NZ000937470	TARA HILLS	-44.517	169.900	TL000091724	NUKUNONO	-9.200	-171.917	4299.297371
NZ000937470	TARA HILLS	-44.517	169.900	NZ000939450	CAMPBELL ISLAND AWS	-52.550	169.167	894.846229
NZ000937470	TARA HILLS	-44.517	169.900	NZM00093439	WELLINGTON AERO AWS	-41.333	174.800	533.227413
NZ000937470	TARA HILLS	-44.517	169.900	NZM00093929	ENDERBY ISLAND AWS	-50.483	166.300	716.176359
NZ000939870	CHATHAM ISLANDS AWS	-43.950	-176.567	TL000091724	NUKUNONO	-9.200	-171.917	3890.098209
NZ000939870	CHATHAM ISLANDS AWS	-43.950	-176.567	NZM00093439	WELLINGTON AERO AWS	-41.333	174.800	763.267727

(b) closest station pair in NZ

```
1 | closest_pair = pairs.orderBy(F.col("DIST_KM").asc()).limit(1)
2 | closest_pair.show(truncate=False)
3 |
4 | stationA = closest_pair.first()["ID_A"]
5 | stationB = closest_pair.first()["ID_B"]
6 | print(f"The closest station pair in New Zealand is station {stationA} and {stationB}")
```

```
1 | +-----+-----+-----+-----+-----+-----+-----+-----+
2 | |ID_A      |NAME_A                |LAT_A|LON_A |ID_B      |NAME_B                |LAT_B |LON_B|DIST_KM      |
3 | +-----+-----+-----+-----+-----+-----+-----+-----+
4 | |NZ000093417|PARAPARAMU AWS        |-40.9|174.983|NZM00093439|WELLINGTON AERO AWS    |-41.333|174.8|50.52909627580285|
5 | +-----+-----+-----+-----+-----+-----+-----+-----+
6 |
7 | The closest station pair in New Zealand is station NZ000093417 and NZM00093439
```

Q3 Daily Climate Summaries Study

(a) Core Element stat

```
1 | # load daily
2 | daily_schema = StructType(
3 |     [
4 |         StructField("ID", StringType()), # Character Station code
5 |         StructField(
6 |             "DATE", StringType()
7 |         ), # Date Observation date formatted as YYYYMMDD
8 |         StructField("ELEMENT", StringType()), # Character Element type indicator
```

```
9 StructField("VALUE", DoubleType()), # Real Data value for ELEMENT
10 StructField("MEASUREMENT", StringType()), # Character Measurement Flag
11 StructField("QUALITY", StringType()), # Character Quality Flag
12 StructField("SOURCE", StringType()), # Character Source Flag
13 StructField("TIME", StringType()), # Time Observation time formatted as HHMM
14 ]
15 )
16
17 daily = spark.read.csv(paths["daily"], schema=daily_schema)
18
19 show_as_html(daily)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

ID	DATE	ELEMENT	VALUE	MEASUREMENT	QUALITY	SOURCE	TIME
ASN00030019	20100101	PRCP	24.0	None	None	a	None
ASN00030021	20100101	PRCP	200.0	None	None	a	None
ASN00030022	20100101	TMAX	294.0	None	None	a	None
ASN00030022	20100101	TMIN	215.0	None	None	a	None
ASN00030022	20100101	PRCP	408.0	None	None	a	None
ASN00029121	20100101	PRCP	820.0	None	None	a	None
ASN00029126	20100101	TMAX	371.0	None	None	S	None
ASN00029126	20100101	TMIN	225.0	None	None	S	None
ASN00029126	20100101	PRCP	0.0	None	None	a	None
ASN00029126	20100101	TAVG	298.0	H	None	S	None
ASN00029127	20100101	TMAX	371.0	None	None	a	None
ASN00029127	20100101	TMIN	225.0	None	None	a	None
ASN00029127	20100101	PRCP	8.0	None	None	a	None
ASN00029129	20100101	PRCP	174.0	None	None	a	None
ASN00029130	20100101	PRCP	86.0	None	None	a	None
ASN00029131	20100101	PRCP	56.0	None	None	a	None
ASN00029132	20100101	PRCP	800.0	None	None	a	None
ASN00029136	20100101	PRCP	22.0	None	None	a	None
ASN00029137	20100101	PRCP	0.0	None	None	a	None
ASN00029139	20100101	TMAX	298.0	None	None	a	None

```
1 # filter daily records containing only the five core elements
2 core_elems = ["TMAX", "TMIN", "PRCP", "SNOW", "SNWD"]
3 daily_core = daily.filter(F.col("ELEMENT").isin(core_elems))
4 show_as_html(daily_core)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

ID	DATE	ELEMENT	VALUE	MEASUREMENT	QUALITY	SOURCE	TIME
ASN00030019	20100101	PRCP	24.0	None	None	a	None



ASN00030021	20100101	PRCP	200.0	None	None	a	None
ASN00030022	20100101	TMAX	294.0	None	None	a	None
ASN00030022	20100101	TMIN	215.0	None	None	a	None
ASN00030022	20100101	PRCP	408.0	None	None	a	None
ASN00029121	20100101	PRCP	820.0	None	None	a	None
ASN00029126	20100101	TMAX	371.0	None	None	S	None
ASN00029126	20100101	TMIN	225.0	None	None	S	None
ASN00029126	20100101	PRCP	0.0	None	None	a	None
ASN00029127	20100101	TMAX	371.0	None	None	a	None
ASN00029127	20100101	TMIN	225.0	None	None	a	None
ASN00029127	20100101	PRCP	8.0	None	None	a	None
ASN00029129	20100101	PRCP	174.0	None	None	a	None
ASN00029130	20100101	PRCP	86.0	None	None	a	None
ASN00029131	20100101	PRCP	56.0	None	None	a	None
ASN00029132	20100101	PRCP	800.0	None	None	a	None
ASN00029136	20100101	PRCP	22.0	None	None	a	None
ASN00029137	20100101	PRCP	0.0	None	None	a	None
ASN00029139	20100101	TMAX	298.0	None	None	a	None
ASN00029139	20100101	TMIN	270.0	None	None	a	None

```

1  # groupby ELEMENT count
2  elem_counts = (
3      daily_core.groupby("ELEMENT")
4      .agg(F.count("*").alias("OBSERVATION_COUNT"))
5      .orderBy(F.desc("OBSERVATION_COUNT"))
6  )
7
8  show_as_html(elem_counts)

```

```

1  .dataframe tbody tr th {
2      vertical-align: top;
3  }
4
5  .dataframe thead th {
6      text-align: right;
7  }

```

ELEMENT	OBSERVATION_COUNT
PRCP	1084610240
TMAX	461915395
TMIN	460752965
SNOW	361688529
SNWD	302055219

```

1  print(f"The most observations is PRCP.")

```

```

1  The most observations is PRCP.

```

(b) TMIN TMAX

```

1 # filter records which reports TMAX but TMIN: base on (ID, DATE)
2 tmax = daily.filter(F.col("ELEMENT") == "TMAX").select(
3     F.col("ID").alias("ID_T"), F.col("DATE").alias("DATE_T")
4 )
5
6 tmin = daily.filter(F.col("ELEMENT") == "TMIN").select(
7     F.col("ID").alias("ID_N"), F.col("DATE").alias("DATE_N")
8 )

```

```

1 # left_anti join to find the part which left(tmax) have but right(tmin) haven't
2 tmax_no_tmin = tmax.join(
3     tmin, (tmax.ID_T == tmin.ID_N) & (tmax.DATE_T == tmin.DATE_N), how="left_anti"
4 )
5 tmax_no_tmin_count = tmax_no_tmin.count()
6
7 show_as_html(tmax_no_tmin)

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	ID_T	DATE_T
AE000041196		19560704
AE000041196		19570607
AE000041196		19580617
AE000041196		19590525
AE000041196		19650714
AE000041196		19650813
AE000041196		19770612
AE000041196		19770904
AE000041196		19780117
AE000041196		19780619
AE000041196		19780910
AE000041196		19780918
AE000041196		19781018
AE000041196		19781024
AE000041196		19781113
AE000041196		19790215
AE000041196		19790420
AE000041196		19790508
AE000041196		19790517
AE000041196		19790609

```

1 print(f"Number of TMAX observations without corresponding TMIN: {tmax_no_tmin_count}")

```

```

1 Number of TMAX observations without corresponding TMIN: 10735252

```

```

1 # unique stations which contribute to missing_pairs count
2 unique_stations_missing = tmax_no_tmin.select("ID_T").distinct().count()
3 show_as_html(unique_stations_missing)
4
5 print(f"tmax_no_tmin_count = {tmax_no_tmin_count}")
6 print(f"unique_stations_missing = {unique_stations_missing}")

```

```

1 [Stage 233:==>          (16 + 8) / 107][Stage 234:>          (0 + 0) / 107]

```

```

1 print(f"Number of unique stations contributing to TMAX observations without corresponding TMIN: {unique_stations_missing}")

```

```

1 Number of unique stations contributing to TMAX observations without corresponding TMIN: 28751

```

```

1 stop_spark()

```

```

1 25/09/12 01:19:46 WARN ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed.

```

## Spark

The spark session is **stopped**, confirm that `yxi75 (notebook)` is under the completed applications section in the Spark UI.

- [Spark UI](#)

```

1

```