

Spark notebook

This notebook will only work in a Jupyter notebook or Jupyter lab session running on the cluster master node in the cloud.

Follow the instructions on the computing resources page to start a cluster and open this notebook.

Steps

1. Connect to the Windows server using Windows App.
2. Connect to Kubernetes.
3. Start Jupyter and open this notebook from Jupyter in order to connect to Spark.

```
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>
```

```

70         f"<ul>",
71         f'<li><a href="http://localhost:{sc.uiWebUrl.split(":")[-1]}" target="_blank">Spark Application UI</a>',
72         f"</ul>",
73         f"<p><b>Config</b></p>",
74         dict_to_html(dict(sc.getConf().getAll())),
75         f"<p><b>Notes</b></p>",
76         f"<ul>",
77         f'<li>The spark session <code>spark</code> and spark context <code>sc</code> global variables have bee',
78         f'<li>Please run <code>stop_spark()</code> before closing the notebook or restarting the kernel or kil',
79         f"</ul>",
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 # Functions to start and stop spark
95
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 17:59:21 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-jav
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.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.startTime	1757570361973
spark.submit.deployMode	client
spark.master	k8s://https://kubernetes.default.svc.cluster.local:443
spark.app.id	spark-7adcb5083a5b42329c59ab38a40282f7
	-Djava.net.preferIPv6Addresses=false --XX:+IgnoreUnrecognizedVMOptions --add-opens=java.base/java.lang=ALL-UNNAMED --add-

spark.driver.extraJavaOptions	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.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.app.submitTime	1757570361858
spark.kubernetes.memoryOverheadFactor	0.3
spark.driver.host	spark-master-svc
spark.kubernetes.executor.podNamePrefix	yxi75-notebook-78045b99375b6dae
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=GzRlhq7EJ0lRHj92oD0lMBNjkc602nrpfb5H8C17FFY%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"

```

```

1 # Use the hdfs command to explore the data in Azure Blob Storage
2 # hdfs dfs-ls -h :Format file sizes in a human-readable fashion (eg 64.0m instead of 67108864).
3 # reference: https://hadoop.apache.org/docs/r3.3.5/hadoop-project-dist/hadoop-common/FileSystemShell.html#ls
4
5
6 # show DATA_ROOT Structure
7 !hdfs dfs -ls -h {DATA_ROOT}

```

```

1 Found 5 items
2 drwxrwxrwx - 0 1970-01-01 12:00 wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd/daily
3 -rwxrwxrwx 1 3.6 K 2025-08-01 21:31 wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd/ghcnd-count
4 -rwxrwxrwx 1 33.6 M 2025-08-01 21:31 wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd/ghcnd-inver
5 -rwxrwxrwx 1 1.1 K 2025-08-01 21:31 wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd/ghcnd-state
6 -rwxrwxrwx 1 10.6 M 2025-08-01 21:31 wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd/ghcnd-stati

```

```

1 # show Daily folder structure and export to txt file
2 !hdfs dfs -ls {paths["daily"]} > ./supplementary/daily_folder_index.txt

```

```

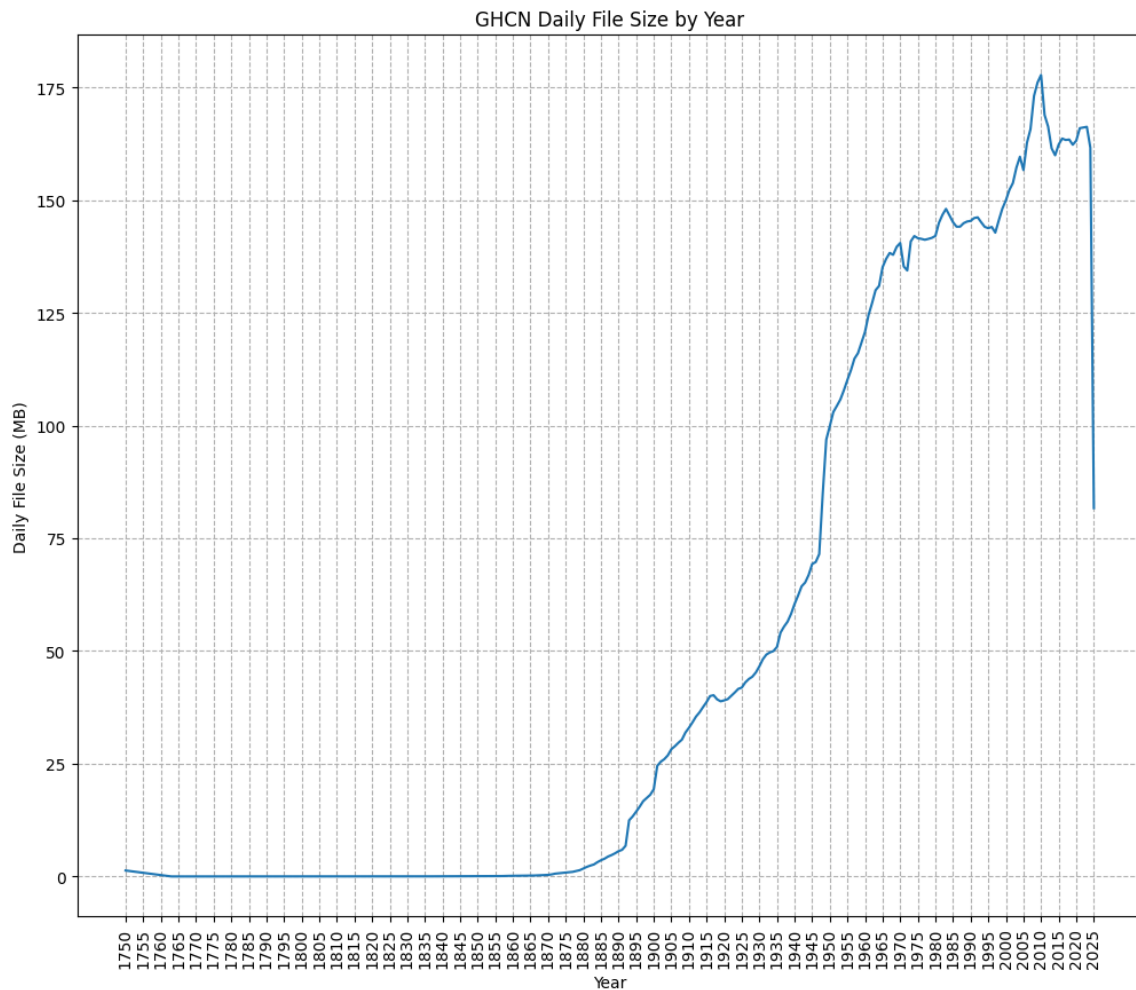
1 # load daily folder index info and plot to show data scale trend
2 daily_folder_index_path = "./supplementary/daily_folder_index.txt"
3 with open(file=daily_folder_index_path, mode="r") as f:
4     lines = f.readlines()
5
6 data = []
7
8 for line in lines:
9     # Skip "Found" line and directory entries
10    if line.startswith("Found") or line.startswith("d"):
11        continue
12
13    # Unpack line into variables
14    file_permissions, replication, size_bytes, mod_date, mod_time, path = line.split()
15
16    # Append row dictionary
17    data.append(
18        {
19            "file_permissions": file_permissions,
20            "replication": int(replication),
21            "size_bytes": int(size_bytes),
22            "mod_date": mod_date,
23            "mod_time": mod_time,
24            "path": path,
25        }
26    )
27
28 # Convert list of dicts to DataFrame

```

```

29 df = pd.DataFrame(data)
30
31 # Add size in MB and extract year
32 df["size_MB"] = df["size_bytes"] / (1024 * 1024)
33 df["year"] = df.path.str.split("/").str[5].str.split(".").str[0].astype(int)
34
35 import matplotlib.pyplot as plt
36
37 plt.figure(figsize=(12, 10))
38 plt.plot(df["year"], df["size_MB"])
39 plt.xticks(range(df["year"].min(), df["year"].max() + 1, 5))
40 plt.xticks(rotation=90)
41 plt.xlabel("Year")
42 ymin = df["size_MB"].min()
43 ymax = df["size_MB"].max()
44 plt.ylabel("Daily File Size (MB)")
45 plt.title("GHCN Daily File Size by Year")
46 plt.grid(visible=True, axis="both", linestyle="--")
47 plt.savefig("../supplementary/GHCN Daily File Size by Year.png", dpi=300)
48 plt.show()

```



Processing

Q1 Data loading and simple EDA

First you will investigate the daily, stations, states, countries, and inventory data provided in cloud storage in [wasbs://campus-data@madstorage002.blob.core.windows.net/ghcnd/](https://campus-data@madstorage002.blob.core.windows.net/ghcnd/) using the `hdfs` command.

Follow the instructions in the notebook provided to explore each dataset using the `hdfs` command without loading any data into memory and answer the following questions:

(a) How is the data structured? Are any of the datasets compressed?

Found 5 items under the directory '[wasbs://campus-data@madstorage002.blob.core.windows.net/ghcnd/](https://campus-data@madstorage002.blob.core.windows.net/ghcnd/)', including 1 directory and 4 text files as shown below:

- daily (264 csv.gz compressed csv files)

- 1750.csv.gz
- 1751.csv.gz
- ...
- 2025.csv.gz
- ghcnd-countries.txt
- ghcnd-inventory.txt
- ghcnd-states.txt
- ghcnd-stations.txt

All the CSV files in the 'daily' directory are indeed compressed using the .gz format.

(b) How many years are contained in daily, and how does the size of the data change?

From the year 1750 to the year 2025, there should be 276 years, but only 264 csv.gz files were found, missing daa for 12 years(1751-1762). Generally, the size of the data increased yearly except for few years, from 10^{-6} MB to 10^5 MB. By the end of 2024, the size of the data reached to 10^5 MB scale.

```
1 # find missing years
2 all_years = set(range(df["year"].min(), df["year"].max() + 1))
3 print("all year count:", len(all_years))
4
5 exisisting_years = set(df["year"])
6 print("GHCN daily dataset year count:", len(exisisting_years))
7
8 missing_year = sorted(all_years - exisisting_years)
9 print("Missing Years:", missing_year)
10 print("Missing Years Count:", len(missing_year))
```

```
1 all year count: 276
2 GHCN daily dataset year count: 264
3 Missing Years: [1751, 1752, 1753, 1754, 1755, 1756, 1757, 1758, 1759, 1760, 1761, 1762]
4 Missing Years Count: 12
```

(c) What is the total size of the data, and how much of that is daily?

```
1 """
2 Usage: hadoop fs -du [-s] [-h] [-v] [-x] URI [URI ...]
3 Displays sizes of files and directories contained in the given directory or the length of a file in case its just a
4 Options:
5 The -s option will result in an aggregate summary of file lengths being displayed, rather than the individual files.
6 The -h option will format file sizes in a "human-readable" fashion (e.g 64.0m instead of 67108864)
7 The -v option will display the names of columns as a header line.
8 The -x option will exclude snapshots from the result calculation. Without the -x option (default), the result is al
9 """
10
11 !hdfs dfs -du -s -h -v {DATA_ROOT}
12 !hdfs dfs -du -s -h -v {paths["daily"]}
```

1	SIZE	DISK_SPACE_CONSUMED_WITH_ALL_REPLICAS	FULL_PATH_NAME
2	13.1 G	13.1 G	wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd
3	SIZE	DISK_SPACE_CONSUMED_WITH_ALL_REPLICAS	FULL_PATH_NAME
4	13.0 G	13.0 G	wasbs://campus-data@mdsstorage002.blob.core.windows.net/ghcnd/daily

Q2 DataType and Schema

You will now load each dataset to ensure the descriptions are accurate and that you can apply the schema either as the data is loaded or by casting columns as they are extracted by manually processing the text records. Extend the code example in the notebook provided by following the steps below.

(a) Define a schema for daily based on the description above or in the GHCN Daily README, using the types defined in pyspark.sql. What do you think is the best way to load the DATE and OBSERVATION TIME columns?

```
1 # Define daily_schema
2 # https://spark.apache.org/docs/latest/api/python/reference/pyspark.sql/api/pyspark.sql.types.StructField.html
3 daily_schema = StructType(
4     [
5         StructField("ID", StringType()), # Character Station code
6         StructField(
7             "DATE", StringType()
8         ), # Date Observation date formatted as YYYYMMDD
9         StructField("ELEMENT", StringType()), # Character Element type indicator
10        StructField("VALUE", DoubleType()), # Real Data value for ELEMENT
```

```

11     StructField("MEASUREMENT", StringType()), # Character Measurement Flag
12     StructField("QUALITY", StringType()), # Character Quality Flag
13     StructField("SOURCE", StringType()), # Character Source Flag
14     StructField("TIME", StringType()), # Time Observation time formatted as HHMM
15 ]
16 )
17
18 # load daily
19 daily = spark.read.csv(
20     path=paths["daily"],
21     schema=daily_schema,
22 )
23
24 print(type(daily))
25 daily.printSchema()
26 daily.show(20)

```

```

1 <class 'pyspark.sql.dataframe.DataFrame'>
2 root
3 |-- ID: string (nullable = true)
4 |-- DATE: string (nullable = true)
5 |-- ELEMENT: string (nullable = true)
6 |-- VALUE: double (nullable = true)
7 |-- MEASUREMENT: string (nullable = true)
8 |-- QUALITY: string (nullable = true)
9 |-- SOURCE: string (nullable = true)
10 |-- TIME: string (nullable = true)
11
12 +-----+-----+-----+-----+-----+-----+-----+-----+
13 |      ID|      DATE|ELEMENT|VALUE|MEASUREMENT|QUALITY|SOURCE|TIME|
14 +-----+-----+-----+-----+-----+-----+-----+-----+
15 |ASN00030019|20100101|PRCP|24.0|      NULL|      NULL|a|NULL|
16 |ASN00030021|20100101|PRCP|200.0|      NULL|      NULL|a|NULL|
17 |ASN00030022|20100101|TMAX|294.0|      NULL|      NULL|a|NULL|
18 |ASN00030022|20100101|TMIN|215.0|      NULL|      NULL|a|NULL|
19 |ASN00030022|20100101|PRCP|408.0|      NULL|      NULL|a|NULL|
20 |ASN00029121|20100101|PRCP|820.0|      NULL|      NULL|a|NULL|
21 |ASN00029126|20100101|TMAX|371.0|      NULL|      NULL|S|NULL|
22 |ASN00029126|20100101|TMIN|225.0|      NULL|      NULL|S|NULL|
23 |ASN00029126|20100101|PRCP| 0.0|      NULL|      NULL|a|NULL|
24 |ASN00029126|20100101|TAVG|298.0|      H|      NULL|S|NULL|
25 |ASN00029127|20100101|TMAX|371.0|      NULL|      NULL|a|NULL|
26 |ASN00029127|20100101|TMIN|225.0|      NULL|      NULL|a|NULL|
27 |ASN00029127|20100101|PRCP| 8.0|      NULL|      NULL|a|NULL|
28 |ASN00029129|20100101|PRCP|174.0|      NULL|      NULL|a|NULL|
29 |ASN00029130|20100101|PRCP| 86.0|      NULL|      NULL|a|NULL|
30 |ASN00029131|20100101|PRCP| 56.0|      NULL|      NULL|a|NULL|
31 |ASN00029132|20100101|PRCP|800.0|      NULL|      NULL|a|NULL|
32 |ASN00029136|20100101|PRCP| 22.0|      NULL|      NULL|a|NULL|
33 |ASN00029137|20100101|PRCP| 0.0|      NULL|      NULL|a|NULL|
34 |ASN00029139|20100101|TMAX|298.0|      NULL|      NULL|a|NULL|
35 +-----+-----+-----+-----+-----+-----+-----+-----+
36 only showing top 20 rows

```

(b) Modify the `spark.read.csv` command to load a subset of the most recent year of daily into Spark so that it uses the schema that you defined in step (a). Did anything go wrong when you tried to use the schema? What data types did you end up using and why?

```

1 latest_year = 2025
2 print(f"load year {latest_year} of daily and show its data schema")
3
4 daily_2025 = spark.read.csv(
5     path=f"{DATA_ROOT}/daily/{latest_year}.csv.gz", schema=daily_schema
6 )
7
8 # parse DATE col
9 daily_2025 = daily_2025.withColumn("DATE", F.to_date(F.col("DATE"), "yyyyMMdd"))
10
11 # parse TIME
12 # 1) DATE: YYYYMMDD -> DateType
13 daily_2025 = daily_2025.withColumn("DATE", F.to_date(F.col("DATE"), "yyyyMMdd"))
14
15 # 2) process TIME col: Null/space/4bit, clean and lpad to 4bit
16 daily_2025 = daily_2025.withColumn(
17     "TIME_CLEAN",
18     F.when(F.col("TIME").isNull() | (F.trim(F.col("TIME")) == ""), None)
19     .otherwise(F.lpad(F.col("TIME").cast("string"), 4, "0"))
20 )

```



```

21
22 # 3) combine "virture date 19700101 + HHmm" and convert to TimestampType
23 daily_2025 = daily_2025.withColumn(
24     "OBS_TS",
25     F.when(F.col("TIME_CLEAN").isNull(), None)
26     .otherwise(F.to_timestamp(
27         F.concat(F.lit("1970-01-01 "), F.col("TIME_CLEAN")), # -> "1970-01-01 HHmm"
28         "yyyy-MM-dd HHmm"
29     ))
30 )
31
32 # 4) "HH:mm" string col for show
33 daily_2025 = daily_2025.withColumn(
34     "OBS_HHMM",
35     F.when(F.col("OBS_TS").isNull(), None)
36     .otherwise(F.date_format(F.col("OBS_TS"), "HH:mm"))
37 )
38
39 # show the sample of data and print the schema
40 print(f"type of daily_latest variable:", type(daily_2025))
41 daily_2025.printSchema()
42 show_as_html(daily_2025)

```

```

1 load year 2025 of daily and show its data schema
2 type of daily_latest variable: <class 'pyspark.sql.dataframe.DataFrame'>
3 root
4 |-- ID: string (nullable = true)
5 |-- DATE: date (nullable = true)
6 |-- ELEMENT: string (nullable = true)
7 |-- VALUE: double (nullable = true)
8 |-- MEASUREMENT: string (nullable = true)
9 |-- QUALITY: string (nullable = true)
10 |-- SOURCE: string (nullable = true)
11 |-- TIME: string (nullable = true)
12 |-- TIME_CLEAN: string (nullable = true)
13 |-- OBS_TS: timestamp (nullable = true)
14 |-- OBS_HHMM: string (nullable = true)

```

```

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	TIME_CLEAN	OBS_TS	OBS_HH
ASN00030019	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00030021	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00030022	2025-01-01	TMAX	414.0	None	None	a	None	None	NaT	None
ASN00030022	2025-01-01	TMIN	247.0	None	None	a	None	None	NaT	None
ASN00030022	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00030025	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029118	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029121	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029126	2025-01-01	TMAX	414.0	None	None	S	None	None	NaT	None
ASN00029126	2025-01-01	TMIN	198.0	None	None	S	None	None	NaT	None
ASN00029126	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029126	2025-01-01	TAVG	321.0	H	None	S	None	None	NaT	None
ASN00029127	2025-01-01	TMAX	414.0	None	None	a	None	None	NaT	None
ASN00029127	2025-01-01	TMIN	198.0	None	None	a	None	None	NaT	None
ASN00029127	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029129	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029131	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029132	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029136	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029139	2025-01-01	TMAX	347.0	None	None	a	None	None	NaT	None

(c) Load each of stations, states, countries, and inventory datasets into Spark and find a way to extract the columns and data types in the descriptions above. You will need to parse the fixed width text formatting by hand, as there is no method to load this format implemented in the standard spark.read library. You should use `pyspark.sql.functions.substring` to extract the columns based on their character range.

```

1  # Q2c: parse fixed-width text (stations/countries/states/inventory)
2  """
3  IV. FORMAT OF "ghcnd-stations.txt"
4  -----
5  Variable   Columns   Type
6  -----
7  ID         1-11     Character
8  LATITUDE   13-20     Real
9  LONGITUDE  22-30     Real
10 ELEVATION  32-37     Real
11 STATE      39-40     Character
12 NAME       42-71     Character
13 GSN FLAG   73-75     Character
14 HCN/CRN FLAG 77-79     Character
15 WMO ID     81-85     Character
16 -----
17 """
18
19 stations_raw = spark.read.text(paths["stations"])
20 stations_df = (
21     stations_raw.withColumn("ID", F.substring("value", 1, 11))
22     .withColumn("LATITUDE", F.substring("value", 13, 8).cast("double"))
23     .withColumn("LONGITUDE", F.substring("value", 22, 9).cast("double"))
24     .withColumn("ELEVATION", F.substring("value", 32, 6).cast("double"))
25     .withColumn("STATE", F.substring("value", 39, 2))
26     .withColumn("NAME", F.substring("value", 42, 30))
27     .withColumn("GSN_FLAG", F.substring("value", 73, 3))
28     .withColumn("HCN_CRN", F.substring("value", 77, 3))
29     .withColumn("WMO_ID", F.substring("value", 81, 5))
30     .drop("value")
31 )
32
33 show_as_html(stations_df)

```

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

ID	LATITUDE	LONGITUDE	ELEVATION	STATE	NAME	GSN_FLAG	HCN_CRN	WMO_ID
ACW00011604	17.1167	-61.7833	10.1		ST JOHNS COOLIDGE FLD			
ACW00011647	17.1333	-61.7833	19.2		ST JOHNS			
AE000041196	25.3330	55.5170	34.0		SHARJAH INTER. AIRP	GSN		41196
AEM00041194	25.2550	55.3640	10.4		DUBAI INTL			41194
AEM00041217	24.4330	54.6510	26.8		ABU DHABI INTL			41217
AEM00041218	24.2620	55.6090	264.9		AL AIN INTL			41218
AF000040930	35.3170	69.0170	3366.0		NORTH-SALANG	GSN		40930
AFM00040938	34.2100	62.2280	977.2		HERAT			40938
AFM00040948	34.5660	69.2120	1791.3		KABUL INTL			40948
AFM00040990	31.5000	65.8500	1010.0		KANDAHAR AIRPORT			40990
AG000060390	36.7167	3.2500	24.0		ALGER-DAR EL BEIDA	GSN		60390
AG000060590	30.5667	2.8667	397.0		EL-GOLEA	GSN		60590
AG000060611	28.0500	9.6331	561.0		IN-AMENAS	GSN		60611
AG000060680	22.8000	5.4331	1362.0		TAMANRASSET	GSN		60680
AGE00135039	35.7297	0.6500	50.0		ORAN-HOPITAL MILITAIRE			
AGE00147704	36.9700	7.7900	161.0		ANNABA-CAP DE GARDE			
AGE00147705	36.7800	3.0700	59.0		ALGIERS-VILLE/UNIVERSITE			
AGE00147706	36.8000	3.0300	344.0		ALGIERS-BOUZAREAH			
AGE00147707	36.8000	3.0400	38.0		ALGIERS-CAP CAXINE			
AGE00147708	36.7200	4.0500	222.0		TIZI OUZOU			60395

```
1 | """
2 | V. FORMAT OF "ghcnd-countries.txt"
3 |
4 | -----
5 | Variable    Columns    Type
6 | -----
7 | CODE         1-2      Character
8 | NAME         4-64     Character
9 | -----
10 |
11 | These variables have the following definitions:
12 |
13 | CODE         is the FIPS country code of the country where the station is
14 |               located (from FIPS Publication 10-4 at
15 |               www.cia.gov/cia/publications/factbook/appendix/appendix-d.html).
16 |
17 | NAME         is the name of the country.
18 | """
19 | countries_df = (
20 |     spark.read.text(paths["countries"])
21 |     .withColumn("CODE", F.substring("value", 1, 2))
22 |     .withColumn("COUNTRY_NAME", F.substring("value", 4, 61))
23 |     .drop("value")
24 | )
25 |
26 | show_as_html(countries_df)
```

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

CODE	COUNTRY_NAME
AC	Antigua and Barbuda
AE	United Arab Emirates
AF	Afghanistan
AG	Algeria
AJ	Azerbaijan
AL	Albania
AM	Armenia
AO	Angola
AQ	American Samoa [United States]
AR	Argentina
AS	Australia
AU	Austria
AY	Antarctica
BA	Bahrain
BB	Barbados
BC	Botswana
BD	Bermuda [United Kingdom]
BE	Belgium
BF	Bahamas, The
BG	Bangladesh

```
1 | """
2 | VI. FORMAT OF "ghcnd-states.txt"
3 |
4 | -----
5 | Variable    Columns    Type
6 | -----
7 | CODE         1-2      Character
8 | NAME         4-50     Character
9 | -----
10 |
11 | These variables have the following definitions:
12 |
13 | CODE         is the POSTAL code of the U.S. state/territory or Canadian
14 |                province where the station is located
15 |
16 | NAME         is the name of the state, territory or province.
17 | """
18 |
19 |
20 | states_df = (
21 |     spark.read.text(paths["states"])
22 |     .withColumn("CODE", F.substring("value", 1, 2))
23 |     .withColumn("STATE_NAME", F.substring("value", 4, 47))
24 |     .drop("value")
25 | )
26 | show_as_html(states_df)
```

```
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
AR	ARKANSAS
AS	AMERICAN SAMOA
AZ	ARIZONA
BC	BRITISH COLUMBIA
CA	CALIFORNIA
CO	COLORADO
CT	CONNECTICUT
DC	DISTRICT OF COLUMBIA
DE	DELAWARE
FL	FLORIDA
FM	MICRONESIA
GA	GEORGIA
GU	GUAM
HI	HAWAII
IA	IOWA
ID	IDAHO
IL	ILLINOIS

```
1 | """
2 | VII. FORMAT OF "ghcnd-inventory.txt"
3 |
4 | -----
5 | Variable    Columns    Type
6 | -----
7 | ID          1-11      Character
8 | LATITUDE    13-20      Real
9 | LONGITUDE   22-30      Real
10 | ELEMENT     32-35      Character
11 | FIRSTYEAR   37-40      Integer
12 | LASTYEAR    42-45      Integer
13 | -----
14 |
15 | These variables have the following definitions:
16 |
17 | ID          is the station identification code. Please see "ghcnd-stations.txt"
18 |              for a complete list of stations and their metadata.
19 |
20 | LATITUDE    is the latitude of the station (in decimal degrees).
21 |
22 | LONGITUDE   is the longitude of the station (in decimal degrees).
23 |
24 | ELEMENT     is the element type. See section III for a definition of elements.
25 |
26 | FIRSTYEAR   is the first year of unflagged data for the given element.
27 |
28 | LASTYEAR    is the last year of unflagged data for the given element.
29 | """
30 |
31 | inventory_df = (
32 |     spark.read.text(paths["inventory"])
```

```

33     .withColumn("ID", F.substring("value", 1, 11))
34     .withColumn("LATITUDE", F.substring("value", 13, 8).cast("double"))
35     .withColumn("LONGITUDE", F.substring("value", 22, 9).cast("double"))
36     .withColumn("ELEMENT", F.substring("value", 32, 4))
37     .withColumn("FIRSTYEAR", F.substring("value", 37, 4).cast("int"))
38     .withColumn("LASTYEAR", F.substring("value", 42, 4).cast("int"))
39     .drop("value")
40 )
41 show_as_html(inventory_df)

```

```

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

```

ID	LATITUDE	LONGITUDE	ELEMENT	FIRSTYEAR	LASTYEAR
ACW00011604	17.1167	-61.7833	TMAX	1949	1949
ACW00011604	17.1167	-61.7833	TMIN	1949	1949
ACW00011604	17.1167	-61.7833	PRCP	1949	1949
ACW00011604	17.1167	-61.7833	SNOW	1949	1949
ACW00011604	17.1167	-61.7833	SNWD	1949	1949
ACW00011604	17.1167	-61.7833	PGTM	1949	1949
ACW00011604	17.1167	-61.7833	WDFG	1949	1949
ACW00011604	17.1167	-61.7833	WSFG	1949	1949
ACW00011604	17.1167	-61.7833	WT03	1949	1949
ACW00011604	17.1167	-61.7833	WT08	1949	1949
ACW00011604	17.1167	-61.7833	WT16	1949	1949
ACW00011647	17.1333	-61.7833	TMAX	1961	1961
ACW00011647	17.1333	-61.7833	TMIN	1961	1961
ACW00011647	17.1333	-61.7833	PRCP	1957	1970
ACW00011647	17.1333	-61.7833	SNOW	1957	1970
ACW00011647	17.1333	-61.7833	SNWD	1957	1970
ACW00011647	17.1333	-61.7833	WT03	1961	1961
ACW00011647	17.1333	-61.7833	WT16	1961	1966
AE000041196	25.3330	55.5170	TMAX	1944	2025
AE000041196	25.3330	55.5170	TMIN	1944	2025

(d) How many rows are there in each of the metadata tables?

```

1  # Q2d-e: rows count
2  meta_counts = {
3      "stations": stations_df.count(),
4      "countries": countries_df.count(),
5      "states": states_df.count(),
6      "inventory": inventory_df.count(),
7      "daily": daily.count(),
8  }
9
10 meta_counts

```

```

1  {'stations': 129657,
2   'countries': 219,
3   'states': 74,
4   'inventory': 766784,
5   'daily': 3155140380}

```

```

1 # save meta counts to txt file
2 with open("./supplementary/meta_counts.txt",mode= "w") as f:
3     for k,v in meta_counts.items():
4         f.write(f"{k}:{v}\n")

```

(e) How many rows are there in daily?

Note that this will take a while if you are only using 2 executors and 1 core per executor, and that the amount of driver and executor memory should not matter unless you actually try to cache or collect all of daily. You should not try to cache or collect all of daily.

```

1 meta_counts.get("daily")

```

```

1 3155140380

```

Q3 Metadata Tables Combination

Next you will combine relevant information from the metadata tables by joining on station, state, and country to give an enriched stations table that we can use for filtering based on attributes at a station level.

(a) Extract the two character country code from each station code in stations and store the output as a new column using the withColumn method.

```

1 # Q3a: extract station ID
2 stations_enriched = stations_df.withColumn("COUNTRY_CODE", F.substring("ID", 1, 2))
3 show_as_html(stations_enriched,5)

```

```

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

```

ID	LATITUDE	LONGITUDE	ELEVATION	STATE	NAME	GSN_FLAG	HCN_CRN	WMO_ID	COUNTRY_CODE
ACW00011604	17.1167	-61.7833	10.1		ST JOHNS COOLIDGE FLD				AC
ACW00011647	17.1333	-61.7833	19.2		ST JOHNS				AC
AE000041196	25.3330	55.5170	34.0		SHARJAH INTER. AIRP	GSN		41196	AE
AEM00041194	25.2550	55.3640	10.4		DUBAI INTL			41194	AE
AEM00041217	24.4330	54.6510	26.8		ABU DHABI INTL			41217	AE

(b) LEFT JOIN stations with countries using your output from step (a).

```

1 # Q3b: LEFT JOIN countries
2 stations_enriched = stations_enriched.join(
3     countries_df.withColumnRenamed(
4         "CODE", "COUNTRY_CODE"
5     ), # rename countries_df col as the same as station_enriched col name before on = "colname" called
6     on="COUNTRY_CODE",
7     how="left",
8 )
9 show_as_html(stations_enriched,5)

```

```

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

```

COUNTRY_CODE	ID	LATITUDE	LONGITUDE	ELEVATION	STATE	NAME	GSN_FLAG	HCN_CRN	WMO_ID
--------------	----	----------	-----------	-----------	-------	------	----------	---------	--------

AC	ACW00011604	17.1167	-61.7833	10.1		ST JOHNS COOLIDGE FLD			
AC	ACW00011647	17.1333	-61.7833	19.2		ST JOHNS			
AE	AE000041196	25.3330	55.5170	34.0		SHARJAH INTER. AIRP	GSN		41196
AE	AEM00041194	25.2550	55.3640	10.4		DUBAI INTL			41194
AE	AEM00041217	24.4330	54.6510	26.8		ABU DHABI INTL			41217

(c) LEFT JOIN stations and states, allowing for the fact that state codes are only provided for stations in the US.

```
1 show_as_html(states_df,2)
2
3 # check stations_enriched non-empty col "STATE" status
4 stations_enriched.filter(F.trim(F.col("STATE")).isNull() & (F.trim(F.col("STATE"))!="")).show()
5
```

```
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	

```
1 +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
2 |COUNTRY_CODE|      ID|LATITUDE|LONGITUDE|ELEVATION|STATE|      NAME|GSN_FLAG|HCN_CRN|WMO_ID|      COUNTRY_NAME|
3 +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
4 |      AQ|AQ00914000|-14.3167|-170.7667|    408.4|  AS|AASUFOU    ...|      |      |      |American Samoa [U...|
5 |      AQ|AQ00914005|-14.2667|   -170.65|    182.9|  AS|AFONO     ...|      |      |      |American Samoa [U...|
6 |      AQ|AQ00914021|-14.2667|-170.5833|     6.1|  AS|AMOULI TUTUILA ...|      |      |      |American Samoa [U...|
7 |      AQ|AQ00914060|-14.2667|-170.6833|    80.8|  AS|ATUU      ...|      |      |      |American Samoa [U...|
8 |      AQ|AQ00914135|   -14.3|   -170.7|   249.9|  AS|FAGA ALU RSVR ...|      |      |      |American Samoa [U...|
9 |      AQ|AQ00914138|-14.2833|-170.6833|    24.1|  AS|FAGA ALU STREAM ...|      |      |      |American Samoa [U...|
10 |      AQ|AQ00914141|-14.2667|-170.6167|     4.6|  AS|FAGAITUA    ...|      |      |      |American Samoa [U...|
11 |      AQ|AQ00914145|-14.2833|-170.7167|    14.9|  AS|FAGASA TUTUILA ...|      |      |      |American Samoa [U...|
12 |      AQ|AQ00914149|-14.2833|-170.6833|    57.0|  AS|FAGA TOGO   ...|      |      |      |American Samoa [U...|
13 |      AQ|AQ00914188|-14.2167|-168.5333|     6.1|  AS|FALEASAO TAU   ...|      |      |      |American Samoa [U...|
14 |      AQ|AQ00914248|-14.2333|-169.5167|     4.9|  AS|FALEASAO VILLAGE ...|      |      |      |American Samoa [U...|
15 |      AQ|AQ00914397|   -14.35|-170.7833|     6.1|  AS|LEONE      ...|      |      |      |American Samoa [U...|
16 |      AQ|AQ00914424|-14.2333|-169.5167|     6.1|  AS|LUMA VILLAGE ...|      |      |      |American Samoa [U...|
17 |      AQ|AQ00914594|-14.3333|-170.7667|    42.4|  AS|MALAELOA    ...|      |      |      |American Samoa [U...|
18 |      AQ|AQ00914650|-14.1667|-169.7167|     3.0|  AS|OFU        ...|      |      |      |American Samoa [U...|
19 |      AQ|AQ00914822|   -11.05|-171.0833|     3.0|  AS|SWAIN ISLAND ...|      |      |      |American Samoa [U...|
20 |      AQ|AQ00914869|-14.3333|-170.7167|     3.0|  AS|TAFUNA AP TUTUILA...|      |      |      |American Samoa [U...|
21 |      AQ|AQ00914873|   -14.35|-170.7667|    14.9|  AS|TAPUTIMU TUTUILA ...|      |      |      |American Samoa [U...|
22 |      AQ|AQ00914902|-14.2728|-170.6922|    80.8|  AS|VAIPITO     ...|      |      |      |American Samoa [U...|
23 |      AQ|AQ00914912|   -14.25|-170.6667|     3.0|  AS|VATIA      ...|      |      |      |American Samoa [U...|
24 +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
25 only showing top 20 rows
```

```
1 # as we already check above cell, we should rename states_df col "CODE" to col "STATE" as to match
2 # station_enriched col "STATE"
3 stations_enriched = stations_enriched.join(
4     states_df.withColumnRenamed("CODE", "STATE"), on="STATE", how="left"
5 )
6
7 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 }
```


STATE	COUNTRY_CODE	ID	LATITUDE	LONGITUDE	ELEVATION	NAME	GSN_FLAG	HCN_CRN	WMO_ID
	AC	ACW00011604	17.1167	-61.7833	10.1	ST JOHNS COOLIDGE FLD			
	AC	ACW00011647	17.1333	-61.7833	19.2	ST JOHNS			
	AE	AE000041196	25.3330	55.5170	34.0	SHARJAH INTER. AIRP	GSN		41196
	AE	AEM00041194	25.2550	55.3640	10.4	DUBAI INTL			41194
	AE	AEM00041217	24.4330	54.6510	26.8	ABU DHABI INTL			41217
	AE	AEM00041218	24.2620	55.6090	264.9	AL AIN INTL			41218
	AF	AF000040930	35.3170	69.0170	3366.0	NORTH-SALANG	GSN		40930
	AF	AFM00040938	34.2100	62.2280	977.2	HERAT			40938
	AF	AFM00040948	34.5660	69.2120	1791.3	KABUL INTL			40948
	AF	AFM00040990	31.5000	65.8500	1010.0	KANDAHAR AIRPORT			40990
	AG	AG000060390	36.7167	3.2500	24.0	ALGER-DAR EL BEIDA	GSN		60390
	AG	AG000060590	30.5667	2.8667	397.0	EL-GOLEA	GSN		60590
	AG	AG000060611	28.0500	9.6331	561.0	IN-AMENAS	GSN		60611
	AG	AG000060680	22.8000	5.4331	1362.0	TAMANRASSET	GSN		60680
	AG	AGE00135039	35.7297	0.6500	50.0	ORAN-HOPITAL MILITAIRE			
	AG	AGE00147704	36.9700	7.7900	161.0	ANNABA-CAP DE GARDE			
	AG	AGE00147705	36.7800	3.0700	59.0	ALGIERS-VILLE/UNIVERSITE			
	AG	AGE00147706	36.8000	3.0300	344.0	ALGIERS-BOUZAREAH			
	AG	AGE00147707	36.8000	3.0400	38.0	ALGIERS-CAP CAXINE			
	AG	AGE00147708	36.7200	4.0500	222.0	TIZI OUZOU			60395

(d) Based on inventory, what was the first and last year that each station was active and collected any element at all?

How many different elements has each station collected overall?
Further, count separately the number of core elements and the number of “other” elements that each station has collected overall.
How many stations collect all five core elements?
How many collect only precipitation and no other elements?

Note that we could also determine the set of elements that each station has collected and store this output as a new column using `pyspark.sql.functions.collect set` but it will be more efficient to first filter inventory by element type using the `element` column and then to join against that output as necessary.

```

1 | show_as_html(inventory_df)

```

```

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

```

ID	LATITUDE	LONGITUDE	ELEMENT	FIRSTYEAR	LASTYEAR
ACW00011604	17.1167	-61.7833	TMAX	1949	1949
ACW00011604	17.1167	-61.7833	TMIN	1949	1949
ACW00011604	17.1167	-61.7833	PRCP	1949	1949
ACW00011604	17.1167	-61.7833	SNOW	1949	1949
ACW00011604	17.1167	-61.7833	SNWD	1949	1949
ACW00011604	17.1167	-61.7833	PGTM	1949	1949
ACW00011604	17.1167	-61.7833	WDFG	1949	1949
ACW00011604	17.1167	-61.7833	WSFG	1949	1949
ACW00011604	17.1167	-61.7833	WT03	1949	1949

ACW00011604	17.1167	-61.7833	WT08	1949	1949
ACW00011604	17.1167	-61.7833	WT16	1949	1949
ACW00011647	17.1333	-61.7833	TMAX	1961	1961
ACW00011647	17.1333	-61.7833	TMIN	1961	1961
ACW00011647	17.1333	-61.7833	PRCP	1957	1970
ACW00011647	17.1333	-61.7833	SNOW	1957	1970
ACW00011647	17.1333	-61.7833	SNWD	1957	1970
ACW00011647	17.1333	-61.7833	WT03	1961	1961
ACW00011647	17.1333	-61.7833	WT16	1961	1966
AE000041196	25.3330	55.5170	TMAX	1944	2025
AE000041196	25.3330	55.5170	TMIN	1944	2025

```

1  # 1) Basic station-level statistics (without using collect_set)
2  base_stats = (
3      inventory_df
4      .groupBy("ID")
5      .agg(
6          F.min("FIRSTYEAR").alias("FIRSTYEAR_ANY"),
7          F.max("LASTYEAR").alias("LASTYEAR_ANY"),
8          F.countDistinct("ELEMENT").alias("N_ELEMENTS")
9      )
10 )
11
12
13 # 2) Filter by element type first (as suggested in the assignment)
14 #   Create one-hot indicator columns for the core elements
15 core = ["TMAX", "TMIN", "PRCP", "SNOW", "SNWD"]
16
17 # Keep only core elements and deduplicate by (station, element)
18 core_pairs = (
19     inventory_df
20     .filter(F.col("ELEMENT").isin(core))
21     .select("ID", "ELEMENT")
22     .distinct()
23 )
24
25 # Pivot to wide format (each core element becomes a 0/1 indicator)
26 # Note: use count(*) > 0 → converted to 1/0
27 core_flags = (
28     core_pairs
29     .groupBy("ID")
30     .pivot("ELEMENT", core)
31     .agg(F.count(F.lit(1)))
32     .na.fill(0)
33 )
34
35 # Optionally: convert counts to explicit 0/1 (can be omitted if already 0/1)
36 for el in core:
37     core_flags = core_flags.withColumn(el, (F.col(el) > 0).cast("int"))
38
39 # -----
40 # 3) Merge: base stats + core element indicators
41 # -----
42 inv_by_station = (
43     base_stats
44     .join(core_flags, on="ID", how="left")
45     .na.fill(0, subset=core) # fill missing core elements with 0
46     .withColumn("N_CORE_ELEMENTS", sum(F.col(el) for el in core))
47 )
48
49 show_as_html(inv_by_station)

```

```

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

```

ID	FIRSTYEAR_ANY	LASTYEAR_ANY	N_ELEMENTS	TMAX	TMIN	PRCP	SNOW	SNWD	N_CORE_ELEMENTS
AE000041196	1944	2025	4	1	1	1	0	0	3
AEM00041218	1994	2025	4	1	1	1	0	0	3
AGE00147715	1879	1938	3	1	1	1	0	0	3
AGE00147717	1880	1938	3	1	1	1	0	0	3
AGE00147719	1888	2025	4	1	1	1	0	0	3
AGM00060425	1943	2025	5	1	1	1	0	1	4
AGM00060430	1957	2025	5	1	1	1	0	1	4
AGM00060490	1957	2025	5	1	1	1	0	1	4
AGM00060514	1995	2025	5	1	1	1	0	1	4
AGM00060515	1984	2025	4	1	1	1	0	0	3
AGM00060522	1976	2025	5	1	1	1	0	1	4
AGM00060536	1981	2025	5	1	1	1	0	1	4
AGM00060555	1958	2025	4	1	1	1	0	0	3
AGM00060580	1957	2025	4	1	1	1	0	0	3
AJ000037742	1955	1987	1	0	0	1	0	0	1
AJ000037747	1955	2025	5	1	1	1	0	1	4
AJ000037756	1955	2025	5	1	1	1	0	1	4
AJ000037816	1963	1991	1	0	0	1	0	0	1
AJ000037899	1936	1991	5	1	1	1	0	1	4
AJ000037914	1955	1991	1	0	0	1	0	0	1

```
1 # # How many stations collect all five core elements?
2 n_all_five_core = inv_by_station.filter(F.col("N_CORE_ELEMENTS") == 5).count()
3
4 n_all_five_core
```

```
1 | 20504
```

```
1 # How many stations collect only precipititation and no other elements?
2 # only PRCP and no other elements
3 n_prcp_only = inv_by_station.filter(
4     (F.col("N_ELEMENTS") == 1) & (F.col("PRCP") == 1)
5 ).count()
6
7 n_prcp_only
```

```
1 | 16267
```

```
1 stations_enriched = stations_enriched.join(
2     inv_by_station.withColumnRenamed("ID", "ID"), on="ID", how="left"
3 )
4
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_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
RSM00030934		RS	50.6000	107.5830	638.0	BICURA			30934
US1VAAP0001	VA	US	37.4083	-78.9638	219.5	CONCORD 4.6 NNE			
US1VAAR0018	VA	US	38.8719	-77.1190	78.0	ARLINGTON 1.0 WSW			
US1VABC0003	VA	US	37.2187	-82.1094	367.3	VANSANT 1.4 SW			
US1VABD0010	VA	US	37.2762	-79.8406	359.7	VINTON 2.6 E			
AG000060590		AG	30.5667	2.8667	397.0	EL-GOLEA	GSN		60590
AGE00147705		AG	36.7800	3.0700	59.0	ALGIERS-VILLE/UNIVERSITE			
AGE00147706		AG	36.8000	3.0300	344.0	ALGIERS-BOUZAREAH			
RSM00030856		RS	51.4830	114.5330	877.0	SEDOVAJA			30856
US1VAAM0002	VA	US	37.3009	-78.1609	118.0	JETERSVILLE 3.6 W			
US1VAAR0005	VA	US	38.8554	-77.0698	41.5	ALEXANDRIA 2.4 NE			
US1VAAR0016	VA	US	38.8920	-77.1617	109.7	ARLINGTON 3.3 WNW			
AF000040930		AF	35.3170	69.0170	3366.0	NORTH-SALANG	GSN		40930
AGE00147704		AG	36.9700	7.7900	161.0	ANNABA-CAP DE GARDE			
AGE00147708		AG	36.7200	4.0500	222.0	TIZI OUZOU			60395
RSM00030862		RS	51.8670	116.0330	597.0	SHILKA			30862
RSM00030935		RS	50.3700	108.7500	770.0	KRASNYJ CHIKOJ			30935
US1VAAR0017	VA	US	38.7986	-77.0431	10.4	HUNTINGTON 1.6 ENE			

20 rows × 21 columns

(e) save combined table stations_enriched in parquet format

```
1 | stations_enriched.count()
2 | size_bytes = stations_enriched.rdd.map(lambda row: len(str(row))).sum()
3 | print("Approx size in MB:", size_bytes / (1024 * 1024))
```

```
1 | [Stage 153:=====> (6 + 3) / 9]
2 |
3 | Approx size in MB: 43.822248458862305
```

```
1 | stations_enriched_savepath = USER_ROOT + "stations_enriched_parquet"
2 | stations_enriched_savepath
3 | # Save the stations_enriched to savepah
4 | stations_enriched.write.mode("overwrite").parquet(stations_enriched_savepath)
```

```
1 | # check save result
2 | !hdfs dfs -ls -h {stations_enriched_savepath}
```

1	Found 10 items									
2	-rw-r--r--	1	yxi75 supergroup	0	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
3	-rw-r--r--	1	yxi75 supergroup	606.8 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
4	-rw-r--r--	1	yxi75 supergroup	582.7 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
5	-rw-r--r--	1	yxi75 supergroup	447.8 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
6	-rw-r--r--	1	yxi75 supergroup	579.3 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
7	-rw-r--r--	1	yxi75 supergroup	459.8 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
8	-rw-r--r--	1	yxi75 supergroup	457.0 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
9	-rw-r--r--	1	yxi75 supergroup	582.2 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
10	-rw-r--r--	1	yxi75 supergroup	589.7 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				
11	-rw-r--r--	1	yxi75 supergroup	452.8 K	2025-09-11 18:30	wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/stations_enriched_				

Q4 Check for missing stations in daily

```
1 # load stations_enriched from Azure storage and cache it for frequently usage
2 stations_enriched = spark.read.parquet(stations_enriched_savepath)
3 stations_enriched.cache()
4 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_CRN	WMO
AE000041196		AE	25.3330	55.5170	34.0	SHARJAH INTER. AIRP	GSN		4119
AEM00041218		AE	24.2620	55.6090	264.9	AL AIN INTL			4121
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			6054
AGM00060425		AG	36.2130	1.3320	141.1	ECH CHELIFF			6042
AGM00060430		AG	36.3000	2.2330	721.0	MILIANA			6043
AGM00060490		AG	35.6240	-0.6210	89.9	ES SENIA			6049
AGM00060514		AG	35.1670	2.3170	801.0	KSAR CHELLALA			6051
AGM00060515		AG	35.3330	4.2060	459.0	BOU SAADA			6051
AGM00060522		AG	34.8200	-1.7700	426.0	MAGHNIA			6052
AGM00060536		AG	34.8670	0.1500	752.0	SAIDA			6053
AGM00060555		AG	33.0680	6.0890	85.0	SIDI MAHDI			6055
AGM00060580		AG	31.9170	5.4130	150.0	OUARGLA			6058
AJ000037742		AJ	40.9000	47.3000	313.0	ORDJONIKIDZE,ZERNOSOVHOZ			3774
AJ000037747		AJ	40.6170	47.1500	15.0	EVLAKH AIRPORT			3774
AJ000037756		AJ	40.5330	48.9330	755.0	MARAZA			3775
AJ000037816		AJ	40.5000	46.1000	1655.0	DASHKESAN			3781
AJ000037899		AJ	39.7670	46.7500	1355.0	SHUSHA			3789
AJ000037914		AJ	39.9000	48.0000	-1.0	IMISLY			3791

20 rows × 21 columns

```
1 # show loaded daily_2025, a subset of daily to check missing id algorithm
2 show_as_html(daily_2025)
```

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

ID	DATE	ELEMENT	VALUE	MEASUREMENT	QUALITY	SOURCE	TIME	TIME_CLEAN	OBS_TS	OBS_HHMM
ASN00030019	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00030021	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00030022	2025-01-01	TMAX	414.0	None	None	a	None	None	NaT	None
ASN00030022	2025-01-01	TMIN	247.0	None	None	a	None	None	NaT	None
ASN00030022	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00030025	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029118	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029121	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029126	2025-01-01	TMAX	414.0	None	None	S	None	None	NaT	None
ASN00029126	2025-01-01	TMIN	198.0	None	None	S	None	None	NaT	None
ASN00029126	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029126	2025-01-01	TAVG	321.0	H	None	S	None	None	NaT	None
ASN00029127	2025-01-01	TMAX	414.0	None	None	a	None	None	NaT	None
ASN00029127	2025-01-01	TMIN	198.0	None	None	a	None	None	NaT	None
ASN00029127	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029129	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029131	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029132	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029136	2025-01-01	PRCP	0.0	None	None	a	None	None	NaT	None
ASN00029139	2025-01-01	TMAX	347.0	None	None	a	None	None	NaT	None

```
1 | # extract year 2025 daily show up station IDs (scalable in yera range)
2 | daily_station_ids_2025 = (
3 |     daily_2025.select("ID").distinct().withColumnRenamed("ID", "ID_IN_DAILY")
4 | )
5 |
6 | # use left anti join to find missing stations in daily_2025
7 | missing_in_daily_2025 = stations_enriched.join(
8 |     daily_station_ids_2025,
9 |     stations_enriched.ID == daily_station_ids_2025.ID_IN_DAILY,
10 |    how="left_anti",
11 | )
12 |
13 | missing_stations_count_2025 = missing_in_daily_2025.count()
14 | missing_stations_count_2025
```

```
1 | 89941
```

```
1 | # extract all the year daily show up station IDs (expand to all yera range)
2 | daily_station_ids_total = (
3 |     daily.select("ID").distinct().withColumnRenamed("ID", "ID_IN_DAILY")
4 | )
5 |
6 | # use left anti join to find missing stations in daily all years
7 | missing_in_daily_total = stations_enriched.join(
8 |     daily_station_ids_total,
```

```
9     stations_enriched.ID == daily_station_ids_total.ID_IN_DAILY,  
10     how="left_anti",  
11 )  
12  
13 missing_stations_count_total = missing_in_daily_total.count()  
14 missing_stations_count_total
```

```
1 | 38
```

```
1 # save to supplementary  
2 with open("./supplementary/missing_in_daily.txt", "w") as f:  
3     f.write(f"missing_stations_count_2025:{missing_stations_count_2025}\n")  
4     f.write(f"missing_stations_count_total:{missing_stations_count_total}\n")  
5
```

```
1 stop_spark()
```

```
1 | 25/09/11 19:26:51 WARN ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed.
```

Spark

The spark session is **stopped**, confirm that `yx175 (notebook)` is under the completed applications section in the Spark UI.

- [Spark UI](#)