

```

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</li>',
77             f'<li>Please run <code>stop_spark()</code> before closing the notebook or restarting the kernel or killing the notebook</li></ul>',
78         ]
79
80     display(HTML("".join(html)))
81

```

```

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=8, executor_cores=2, worker_memory=4, master_memory=6)

```

```

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 25/09/14 12:00:39 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.app.id	spark-77fd2001216d4e54bbde6baa86933163
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.cores.max	16
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.app.submitTime	1757808039790
spark.driver.memory	6g
spark.kubernetes.namespace	yxi75
spark.serializer.objectStreamReset	100
spark.driver.maxResultSize	0
spark.submit.deployMode	client
spark.master	k8s://https://kubernetes.default.svc.cluster.local:443
spark.driver.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 - Dderby.system.home=/tmp/yxi75/spark/
spark.fs.azure	org.apache.hadoop.fs.azure.NativeAzureFileSystem
spark.sql.shuffle.partitions	64
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.instances	8
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=GzRlhq7EJ01RHj92oD01MBNjkc602nrpfB5H8C17FFY%3D
spark.rdd.compress	True
spark.kubernetes.executor.podNamePrefix	yxi75-notebook-4f45329945861b0e
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.driver.port	7077
spark.submit.pyFiles	
spark.app.startTime	1757808039932
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 os
4  import re
5  import subprocess
6  import sys
7  from math import asin, cos, radians, sin, sqrt
8  from pprint import pprint
9
10 import cartopy
11 import cartopy.crs as ccrs

```

```

12 import geopandas as gpd
13 import matplotlib.pyplot as plt
14 import numpy as np
15 import pandas as pd
16 from cartopy.mpl.gridliner import LATITUDE_FORMATTER, LONGITUDE_FORMATTER
17 from pyspark.sql import functions as F
18 from pyspark.sql.functions import udf
19 from pyspark.sql.types import *

```

```

1 # Paths global variables
2 DATA_ROOT = "wasbs://campus-data@madsstorage002.blob.core.windows.net/ghcnd/"
3 USER_ROOT = "wasbs://campus-user@madsstorage002.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 stations_enriched_savepath = USER_ROOT + "stations_enriched_parquet"
14 station_count_by_country_path = USER_ROOT + "station_count_by_country_parquet"
15 station_count_us_terri_path = USER_ROOT + "station_count_us_terri_parquet"
16 country_meta_with_station_num_path = USER_ROOT + "country_meta_with_station_num"
17 states_meta_with_station_num_path = USER_ROOT + "states_meta_with_station_num"
18 daily_nz_tmin_tmax_path = USER_ROOT + "daily_nz_tmin_tmax_parquet"
19 prcp_pdf_path = USER_ROOT + "prcp_pdf_parquet"

```

## plot observation of TMIN and TMAX for stations in NZ

```

1 !hdfs dfs -ls {USER_ROOT}

```

```

1 Found 7 items
2 drwxr-xr-x - yxi75 supergroup          0 2025-09-11 22:55 wasbs://campus-user@madsstorage002.blob.core.windows.net/
3 drwxr-xr-x - yxi75 supergroup          0 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/
4 drwxr-xr-x - yxi75 supergroup          0 2025-09-14 07:29 wasbs://campus-user@madsstorage002.blob.core.windows.net/
5 drwxr-xr-x - yxi75 supergroup          0 2025-09-11 23:18 wasbs://campus-user@madsstorage002.blob.core.windows.net/
6 drwxr-xr-x - yxi75 supergroup          0 2025-09-11 22:33 wasbs://campus-user@madsstorage002.blob.core.windows.net/
7 drwxr-xr-x - yxi75 supergroup          0 2025-09-11 22:33 wasbs://campus-user@madsstorage002.blob.core.windows.net/
8 drwxr-xr-x - yxi75 supergroup          0 2025-09-11 18:30 wasbs://campus-user@madsstorage002.blob.core.windows.net/

```

```

1 station_count_by_country = spark.read.parquet(station_count_by_country_path)

```

```

1 # How many stations we have in NZ
2 station_count_by_country.filter(F.col("COUNTRY_CODE").contains("NZ")).show()
3 # so we only have 15 stations in NZ, not a big number.

```

```

1 +-----+-----+-----+
2 |COUNTRY_CODE|COUNTRY_NAME|count|
3 +-----+-----+-----+
4 |           NZ|New Zealand |   15|
5 +-----+-----+-----+

```

```

1 # now we get NZ stations' ID from stations_enriched first
2 # then inner join with daily dataset to filter out NZ station observations
3 stations_enriched = spark.read.parquet(stations_enriched_savepath)
4
5 nz_station_ids = stations_enriched.filter(F.col("COUNTRY_CODE").contains("NZ")).select(
6     "ID"
7 )

```

```

1 # in order to inner join daily, we should load daily first.
2
3 # Define schma for Daily
4 daily_schema = StructType(
5     [
6         StructField("ID", StringType(), nullable=False),

```

```

7      StructField("DATE", StringType(), nullable=False),
8      StructField("ELEMENT", StringType(), nullable=False),
9      StructField("VALUE", FloatType(), nullable=False),
10     StructField("MEASUREMENT_FLAG", StringType(), nullable=True),
11     StructField("QUALITY_FLAG", StringType(), nullable=True),
12     StructField("SOURCE_FLAG", StringType(), nullable=True),
13     StructField("OBSERVATION_TIME", StringType(), nullable=True),
14   ]
15 )
16
17 # load daily and check daily schema for later join parameter on = ""
18 daily = spark.read.csv(paths["daily"], schema=daily_schema)

```

```

1 # Here's a decision making point
2 # A small table nz_station_ids and a large table daily, when join, choose broadcast join rather than shuffle join.
3 # broadcast nz_station_ids to all partitions for locally join, and pass the join result back to the master node only.
4
5 daily_nz = daily.join(F.broadcast(nz_station_ids), on="ID", how="inner")
6
7 show_as_html(daily_nz)

```

```

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_FLAG	QUALITY_FLAG	SOURCE_FLAG	OBSERVATION_TIME
NZ000093012	20100101	TAVG	178.0	H	None	S	None
NZ000093292	20100101	TMAX	297.0	None	None	S	None
NZ000093292	20100101	TMIN	74.0	None	None	S	None
NZ000093292	20100101	PRCP	0.0	None	None	S	None
NZ000093292	20100101	TAVG	235.0	H	None	S	None
NZ000093417	20100101	TMAX	180.0	None	None	S	None
NZ000093417	20100101	TMIN	125.0	None	None	S	None
NZ000093417	20100101	PRCP	0.0	None	None	S	None
NZ000093417	20100101	TAVG	163.0	H	None	S	None
NZ000093844	20100101	TMAX	232.0	None	None	S	None
NZ000093844	20100101	TMIN	96.0	None	None	S	None
NZ000093844	20100101	PRCP	8.0	None	None	S	None
NZ000093844	20100101	TAVG	167.0	H	None	S	None
NZ000093994	20100101	PRCP	0.0	None	None	S	None
NZ000093994	20100101	TAVG	204.0	H	None	S	None
NZ000933090	20100101	TMAX	197.0	None	None	S	None
NZ000933090	20100101	TMIN	82.0	None	None	S	None
NZ000933090	20100101	PRCP	0.0	None	None	S	None
NZ000933090	20100101	TAVG	168.0	H	None	S	None
NZ000936150	20100101	TMAX	324.0	None	None	S	None

```

1 # filter out TMIN and TMAX from daily_nz, and save it to output path for later plot use
2 daily_nz_tmin_tmax = daily_nz.filter(F.col("ELEMENT").isin(["TMIN", "TMAX"]))
3
4 show_as_html(daily_nz_tmin_tmax)
5
6 daily_nz_tmin_tmax_count = daily_nz_tmin_tmax.count()
7 print(f"daily_nz_tmin_tmax observation count: {daily_nz_tmin_tmax_count}")

```

```
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_FLAG	QUALITY_FLAG	SOURCE_FLAG	OBSERVATION_TIME
NZ000093292	20100101	TMAX	297.0	None	None	S	None
NZ000093292	20100101	TMIN	74.0	None	None	S	None
NZ000093417	20100101	TMAX	180.0	None	None	S	None
NZ000093417	20100101	TMIN	125.0	None	None	S	None
NZ000093844	20100101	TMAX	232.0	None	None	S	None
NZ000093844	20100101	TMIN	96.0	None	None	S	None
NZ000933090	20100101	TMAX	197.0	None	None	S	None
NZ000933090	20100101	TMIN	82.0	None	None	S	None
NZ000936150	20100101	TMAX	324.0	None	None	S	None
NZM00093110	20100101	TMAX	215.0	None	None	S	None
NZM00093110	20100101	TMIN	153.0	None	None	S	None
NZM00093439	20100101	TMAX	204.0	None	None	S	None
NZM00093439	20100101	TMIN	134.0	None	None	S	None
NZM00093678	20100101	TMAX	242.0	None	None	S	None
NZM00093678	20100101	TMIN	94.0	None	None	S	None
NZM00093781	20100101	TMAX	324.0	None	None	S	None
NZ000093292	20100102	TMAX	302.0	None	None	S	None
NZ000093292	20100102	TMIN	180.0	None	None	S	None
NZ000093417	20100102	TMAX	181.0	None	None	S	None
NZ000093417	20100102	TMIN	153.0	None	None	S	None

```
1 | [Stage 10:=====>(106 + 1) / 107]
2 |
3 | daily_nz_tmin_tmax observation count: 494311
```

```
1 | # save to daily_nz_tmin_tmax_path
2 | # daily_nz_tmin_tmax.write.parquet(daily_nz_tmin_tmax_path)
3 |
4 | # check save result
5 | !hdfs dfs -ls -h {daily_nz_tmin_tmax_path}
```

```
1 | Found 86 items
2 | -rw-r--r-- 1 yxi75 supergroup 0 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
3 | -rw-r--r-- 1 yxi75 supergroup 13.7 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
4 | -rw-r--r-- 1 yxi75 supergroup 13.7 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
5 | -rw-r--r-- 1 yxi75 supergroup 13.7 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
6 | -rw-r--r-- 1 yxi75 supergroup 13.8 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
7 | -rw-r--r-- 1 yxi75 supergroup 13.9 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
8 | -rw-r--r-- 1 yxi75 supergroup 16.1 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
9 | -rw-r--r-- 1 yxi75 supergroup 16.0 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
10 | -rw-r--r-- 1 yxi75 supergroup 15.7 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
11 | -rw-r--r-- 1 yxi75 supergroup 13.6 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
12 | -rw-r--r-- 1 yxi75 supergroup 16.5 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
13 | -rw-r--r-- 1 yxi75 supergroup 14.5 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
14 | -rw-r--r-- 1 yxi75 supergroup 15.5 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
15 | -rw-r--r-- 1 yxi75 supergroup 16.2 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
16 | -rw-r--r-- 1 yxi75 supergroup 13.4 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
17 | -rw-r--r-- 1 yxi75 supergroup 14.6 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
18 | -rw-r--r-- 1 yxi75 supergroup 15.2 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
19 | -rw-r--r-- 1 yxi75 supergroup 16.0 K 2025-09-13 13:02 wasbs://campus-user@madsstorage002.blob.core.windows.net/yxi75/daily_nz_tmin_tmax
```





```
1 # how many years are covered by the daily_nz_tmin_tmax
2 # parse DATE col
3 daily_nz_tmin_tmax = daily_nz_tmin_tmax.withColumn(
4     "DATE", F.to_date(F.col("DATE"), "yyyyMMdd")
5 )
6
7 # extract years col
8 years = (
9     daily_nz_tmin_tmax.select(F.year("DATE").alias("year")).distinct().orderBy("year")
10 )
11
12 # calculate year range
13 years_agg = years.agg(
14     F.min("year").alias("min_year"),
15     F.max("year").alias("max_year"),
16     F.countDistinct("year").alias("n_year"),
17 )
18 show_as_html(years_agg)
```

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

	min_year	max_year	n_year
	1940	2025	86

## Check data gap in year level

From 1940 to 2025, there should have  $2025-1940+1=86$  years, where our distinct method found 86 in years\_agg, so no year-level gap.

## Check data gap in month level

```
1 # ===== 0) Set year range =====
2 min_y, max_y = 1940, 2025
3
4 # ===== 1) Dimension tables: station / year / month =====
5 # Station table: you already have nz_station_ids (small table).
6 # If IDs are not guaranteed unique, apply distinct upstream.
7 ids_df = nz_station_ids.select("ID")
8 years_df = spark.range(min_y, max_y + 1).withColumnRenamed("id", "year")
9 months_df = spark.range(1, 13).withColumnRenamed("id", "month")
10
11 # Cartesian product to generate the full (ID, year, month) frame
12 # Expected size = num_stations * num_years * 12 = 15480 row
13 full_frame = ids_df.crossJoin(years_df).crossJoin(months_df)
```

```
1 | show_as_html(full_frame)
```

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

	ID	year	month
	NZ000936150	1940	1
	NZ000936150	1940	2
	NZ000936150	1940	3
	NZ000936150	1940	4
	NZ000936150	1940	5
	NZ000936150	1940	6

NZ000936150	1940	7
NZ000936150	1940	8
NZ000936150	1940	9
NZ000936150	1940	10
NZ000936150	1940	11
NZ000936150	1940	12
NZ000936150	1941	1
NZ000936150	1941	2
NZ000936150	1941	3
NZ000936150	1941	4
NZ000936150	1941	5
NZ000936150	1941	6
NZ000936150	1941	7
NZ000936150	1941	8

```

1  # ===== 2) Observed (ID, year, month) combinations =====
2  # Instead of distinct, aggregate with groupBy once
3  obs_by_month = (
4      daily_nz_tmin_tmax.select(
5          F.col("ID"), F.year("DATE").alias("year"), F.month("DATE").alias("month")
6      )
7      .groupBy("ID", "year", "month")
8      .agg(F.count(F.lit(1)).alias("cnt")) # just to prove existence
9      .select("ID", "year", "month") # reduce back to presence set
10 )

```

```

1  print(f"obs_by_month has {obs_by_month.count()} rows")
2  # obs_by_month is a small table too! only 8954 rows
3
4  obs_by_month.show()

```

```

1  obs_by_month has 8954 rows
2  +-----+-----+
3  |      ID|year|month|
4  +-----+-----+
5  |NZ000093417|2011| 8|
6  |NZM00093439|2011| 10|
7  |NZ000093292|2011| 3|
8  |NZ000093292|2011| 8|
9  |NZ000093417|2011| 11|
10 |NZ000093417|2011| 4|
11 |NZM00093678|2011| 4|
12 |NZM00093678|2011| 6|
13 |NZ000093292|2011| 7|
14 |NZM00093678|2011| 7|
15 |NZ000936150|2011| 10|
16 |NZ000936150|2011| 11|
17 |NZ000093417|2011| 2|
18 |NZ000093844|2011| 7|
19 |NZ000093844|2011| 4|
20 |NZ000093292|2011| 4|
21 |NZ000093417|2011| 5|
22 |NZ000093417|2011| 7|
23 |NZ000093417|2011| 9|
24 |NZM00093678|2011| 10|
25 +-----+-----+
26 only showing top 20 rows

```

```

1 # ===== 3) Find missing (expected minus observed) =====
2 # full_frame is small (15480 rows), use it as left table with left_anti join
3 missing = full_frame.join(obs_by_month, on=["ID", "year", "month"], how="left_anti")
4
5 # ===== 4) Aggregate missing months for each (ID, year) =====
6 gaps_by_id_year = (
7     missing.groupBy("ID", "year").agg(F.collect_list("month").alias("missing_months"))
8 ).orderBy("ID", "year", "missing_months")
9
10 gaps_by_id_year.show(100, truncate=False)

```

```

1 +-----+-----+-----+
2 |ID|year|missing_months|
3 +-----+-----+-----+
4 |NZ000093012|1940|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
5 |NZ000093012|1941|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
6 |NZ000093012|1942|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
7 |NZ000093012|1943|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
8 |NZ000093012|1944|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
9 |NZ000093012|1945|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
10 |NZ000093012|1946|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
11 |NZ000093012|1947|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
12 |NZ000093012|1948|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
13 |NZ000093012|1949|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
14 |NZ000093012|1950|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
15 |NZ000093012|1951|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
16 |NZ000093012|1952|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
17 |NZ000093012|1953|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
18 |NZ000093012|1954|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
19 |NZ000093012|1955|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
20 |NZ000093012|1956|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
21 |NZ000093012|1957|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
22 |NZ000093012|1958|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
23 |NZ000093012|1959|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
24 |NZ000093012|1960|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
25 |NZ000093012|1961|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
26 |NZ000093012|1962|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
27 |NZ000093012|1963|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
28 |NZ000093012|1964|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
29 |NZ000093012|1965|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
30 |NZ000093012|1970|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
31 |NZ000093012|1971|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
32 |NZ000093012|1972|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
33 |NZ000093012|1973|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
34 |NZ000093012|1974|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
35 |NZ000093012|1975|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
36 |NZ000093012|1976|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
37 |NZ000093012|1977|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
38 |NZ000093012|1978|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
39 |NZ000093012|1979|[1, 2, 3]|
40 |NZ000093012|2006|[1, 2, 3, 4, 5, 6, 8, 9, 10, 11, 12]|
41 |NZ000093012|2007|[1, 2, 3, 4, 5, 8, 10, 11, 12]|
42 |NZ000093012|2008|[1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12]|
43 |NZ000093012|2009|[1, 2, 3, 4, 5, 9, 10, 11, 12]|
44 |NZ000093012|2010|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
45 |NZ000093012|2011|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
46 |NZ000093012|2012|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
47 |NZ000093012|2013|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
48 |NZ000093012|2014|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
49 |NZ000093012|2015|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
50 |NZ000093012|2016|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
51 |NZ000093012|2017|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
52 |NZ000093012|2018|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
53 |NZ000093012|2019|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
54 |NZ000093012|2020|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
55 |NZ000093012|2021|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
56 |NZ000093012|2022|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
57 |NZ000093012|2023|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
58 |NZ000093012|2024|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
59 |NZ000093012|2025|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
60 |NZ000093292|1940|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
61 |NZ000093292|1941|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
62 |NZ000093292|1942|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
63 |NZ000093292|1943|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
64 |NZ000093292|1944|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
65 |NZ000093292|1945|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
66 |NZ000093292|1946|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
67 |NZ000093292|1947|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
68 |NZ000093292|1948|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
69 |NZ000093292|1949|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
70 |NZ000093292|1950|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
71 |NZ000093292|1951|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|

```

```
72 | NZ000093292|1952|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
73 | NZ000093292|1953|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
74 | NZ000093292|1954|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
75 | NZ000093292|1955|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
76 | NZ000093292|1956|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
77 | NZ000093292|1957|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
78 | NZ000093292|1958|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
79 | NZ000093292|1959|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
80 | NZ000093292|1960|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
81 | NZ000093292|1961|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
82 | NZ000093292|1962|[1] |
83 | NZ000093292|2025|[8, 9, 10, 11, 12] |
84 | NZ000093417|1940|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
85 | NZ000093417|1941|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
86 | NZ000093417|1942|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
87 | NZ000093417|1943|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
88 | NZ000093417|1944|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
89 | NZ000093417|1945|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
90 | NZ000093417|1946|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
91 | NZ000093417|1947|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
92 | NZ000093417|1948|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
93 | NZ000093417|1949|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
94 | NZ000093417|1950|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
95 | NZ000093417|1951|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
96 | NZ000093417|1952|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
97 | NZ000093417|1953|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
98 | NZ000093417|1954|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
99 | NZ000093417|1955|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
100 | NZ000093417|1956|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
101 | NZ000093417|1957|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
102 | NZ000093417|1958|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
103 | NZ000093417|1959|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
104 | +-----+-----+-----+
105 | only showing top 100 rows
```

## gap months stats

```
1 | # define the latest observation month in our daily dataset, so 2025 Aug or later should not be missing.
2 | obs_by_month.filter(F.col("year") == 2025).agg(F.max(F.col("month"))).show()
```

```
1 | +-----+
2 | |max(month)|
3 | +-----+
4 | |          7|
5 | +-----+
```

```
1 | daily_nz_tmin_tmax_missing_stats = (
2 |   gaps_by_id_year.withColumn(
3 |     "missing_month_count", F.size("missing_months")
4 |   ) # missing in current year
5 |   .withColumn(
6 |     "missing_year_flag", F.when(F.size("missing_months") > 0, 1).otherwise(0)
7 |   ) # if current year has missing month, flag=1
8 |   .groupBy("ID")
9 |   .agg(
10 |     F.lit(F.sum("missing_year_flag")).alias(
11 |       "missing_year_total"
12 |     ), # -5 is because 2025 the latest record is July, so Aug or later is not missing.
13 |     F.lit(F.sum("missing_month_count") - 5).alias("missing_month_total"),
14 |   )
15 | ).orderBy("missing_month_total")
16 |
17 | show_as_html(daily_nz_tmin_tmax_missing_stats)
```

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

	ID	missing_year_total	missing_month_total
	NZ000933090	5	48

NZ000093994	10	87
NZ000093844	10	100
NZ000939450	17	176
NZ000093292	24	265
NZ000936150	25	288
NZ000937470	31	364
NZ000093417	33	384
NZM00093781	40	448
NZ000939870	39	450
NZ000093012	56	650
NZM00093678	64	750
NZM00093110	64	752
NZM00093439	64	752
NZM00093929	81	937

```
1  # verify result by sampling check one station ID
2  gaps_by_id_year.filter(F.col("ID") == "NZ000093844").show(
3    100, truncate=False
4  ) # should missing 100 months
5
6  # check: misisng 1940-1947, 1948 obs has 8 months, matching gaps_by_id_year misisng 4 months in 1948
7  obs_by_month.filter(F.col("ID") == "NZ000093844").groupBy("year").agg(
8    F.count("month")
9  ).orderBy("year").show(100)
```

```
1  +-----+---+-----+-----+
2  |ID          |year|missing_months|
3  +-----+---+-----+-----+
4  |NZ000093844|1940|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
5  |NZ000093844|1941|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
6  |NZ000093844|1942|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
7  |NZ000093844|1943|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
8  |NZ000093844|1944|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
9  |NZ000093844|1945|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
10 |NZ000093844|1946|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
11 |NZ000093844|1947|[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]|
12 |NZ000093844|1948|[1, 2, 3, 4]|
13 |NZ000093844|2025|[8, 9, 10, 11, 12]|
14 +-----+---+-----+-----+
15
16 +---+-----+
17 |year|count(month)|
18 +---+-----+
19 |1948|      8|
20 |1949|     12|
21 |1950|     12|
22 |1951|     12|
23 |1952|     12|
24 |1953|     12|
25 |1954|     12|
26 |1955|     12|
27 |1956|     12|
28 |1957|     12|
29 |1958|     12|
30 |1959|     12|
31 |1960|     12|
32 |1961|     12|
33 |1962|     12|
34 |1963|     12|
35 |1964|     12|
36 |1965|     12|
37 |1966|     12|
38 |1967|     12|
39 |1968|     12|
40 |1969|     12|
41 |1970|     12|
42 |1971|     12|
43 |1972|     12|
44 |1973|     12|
```

45	1974	12
46	1975	12
47	1976	12
48	1977	12
49	1978	12
50	1979	12
51	1980	12
52	1981	12
53	1982	12
54	1983	12
55	1984	12
56	1985	12
57	1986	12
58	1987	12
59	1988	12
60	1989	12
61	1990	12
62	1991	12
63	1992	12
64	1993	12
65	1994	12
66	1995	12
67	1996	12
68	1997	12
69	1998	12
70	1999	12
71	2000	12
72	2001	12
73	2002	12
74	2003	12
75	2004	12
76	2005	12
77	2006	12
78	2007	12
79	2008	12
80	2009	12
81	2010	12
82	2011	12
83	2012	12
84	2013	12
85	2014	12
86	2015	12
87	2016	12
88	2017	12
89	2018	12
90	2019	12
91	2020	12
92	2021	12
93	2022	12
94	2023	12
95	2024	12
96	2025	7
97	+-----+-----+	

```
1 # verify daily_nz_tmin_tmax_missing_stats with daily record sampling
2 daily_nz_tmin_tmax.filter(
3     (F.col("ID") == "NZ000093844") & (F.col("DATE").contains("1948")))
4 ).show(400, truncate=False)
```

1	+-----+-----+-----+-----+-----+-----+-----+-----+							
2	ID	DATE	ELEMENT	VALUE	MEASUREMENT_FLAG	QUALITY_FLAG	SOURCE_FLAG	OBSERVATION_TIME
3	+-----+-----+-----+-----+-----+-----+-----+-----+							
4	NZ000093844	1948-05-31	TMIN	-28.0	NULL	NULL	G	NULL
5	NZ000093844	1948-06-01	TMAX	90.0	NULL	NULL	G	NULL
6	NZ000093844	1948-06-01	TMIN	14.0	NULL	NULL	G	NULL
7	NZ000093844	1948-06-02	TMAX	93.0	NULL	NULL	G	NULL
8	NZ000093844	1948-06-02	TMIN	10.0	NULL	NULL	G	NULL
9	NZ000093844	1948-06-03	TMAX	88.0	NULL	NULL	G	NULL
10	NZ000093844	1948-06-03	TMIN	-32.0	NULL	NULL	G	NULL
11	NZ000093844	1948-06-04	TMAX	94.0	NULL	NULL	G	NULL
12	NZ000093844	1948-06-04	TMIN	12.0	NULL	NULL	G	NULL
13	NZ000093844	1948-06-05	TMAX	84.0	NULL	NULL	G	NULL
14	NZ000093844	1948-06-05	TMIN	17.0	NULL	NULL	G	NULL
15	NZ000093844	1948-06-06	TMAX	100.0	NULL	NULL	G	NULL
16	NZ000093844	1948-06-06	TMIN	-4.0	NULL	NULL	G	NULL
17	NZ000093844	1948-06-07	TMAX	93.0	NULL	NULL	G	NULL
18	NZ000093844	1948-06-07	TMIN	-61.0	NULL	NULL	G	NULL
19	NZ000093844	1948-06-08	TMAX	72.0	NULL	NULL	G	NULL
20	NZ000093844	1948-06-08	TMIN	-22.0	NULL	NULL	G	NULL

21	NZ000093844 1948-06-09 TMAX 92.0 NULL NULL G NULL
22	NZ000093844 1948-06-09 TMIN 45.0 NULL NULL G NULL
23	NZ000093844 1948-06-10 TMAX 106.0 NULL NULL G NULL
24	NZ000093844 1948-06-10 TMIN -2.0 NULL NULL G NULL
25	NZ000093844 1948-06-11 TMAX 114.0 NULL NULL G NULL
26	NZ000093844 1948-06-11 TMIN 8.0 NULL NULL G NULL
27	NZ000093844 1948-06-12 TMAX 117.0 NULL NULL G NULL
28	NZ000093844 1948-06-12 TMIN 38.0 NULL NULL G NULL
29	NZ000093844 1948-06-13 TMAX 89.0 NULL NULL G NULL
30	NZ000093844 1948-06-13 TMIN -39.0 NULL NULL G NULL
31	NZ000093844 1948-06-14 TMAX 88.0 NULL NULL G NULL
32	NZ000093844 1948-06-14 TMIN -44.0 NULL NULL G NULL
33	NZ000093844 1948-06-15 TMAX 106.0 NULL NULL G NULL
34	NZ000093844 1948-06-15 TMIN 22.0 NULL NULL G NULL
35	NZ000093844 1948-06-16 TMAX 97.0 NULL NULL G NULL
36	NZ000093844 1948-06-16 TMIN 2.0 NULL NULL G NULL
37	NZ000093844 1948-06-17 TMAX 100.0 NULL NULL G NULL
38	NZ000093844 1948-06-17 TMIN -9.0 NULL NULL G NULL
39	NZ000093844 1948-06-18 TMAX 94.0 NULL NULL G NULL
40	NZ000093844 1948-06-18 TMIN 6.0 NULL NULL G NULL
41	NZ000093844 1948-06-19 TMAX 61.0 NULL NULL G NULL
42	NZ000093844 1948-06-19 TMIN 17.0 NULL NULL G NULL
43	NZ000093844 1948-06-20 TMAX 87.0 NULL NULL G NULL
44	NZ000093844 1948-06-20 TMIN 7.0 NULL NULL G NULL
45	NZ000093844 1948-06-21 TMAX 89.0 NULL NULL G NULL
46	NZ000093844 1948-06-21 TMIN 47.0 NULL NULL G NULL
47	NZ000093844 1948-06-22 TMAX 107.0 NULL NULL G NULL
48	NZ000093844 1948-06-22 TMIN 5.0 NULL NULL G NULL
49	NZ000093844 1948-06-23 TMAX 128.0 NULL NULL G NULL
50	NZ000093844 1948-06-23 TMIN 16.0 NULL NULL G NULL
51	NZ000093844 1948-06-24 TMAX 83.0 NULL NULL G NULL
52	NZ000093844 1948-06-24 TMIN 32.0 NULL NULL G NULL
53	NZ000093844 1948-06-25 TMAX 79.0 NULL NULL G NULL
54	NZ000093844 1948-06-25 TMIN 36.0 NULL NULL G NULL
55	NZ000093844 1948-06-26 TMAX 81.0 NULL NULL G NULL
56	NZ000093844 1948-06-26 TMIN -27.0 NULL NULL G NULL
57	NZ000093844 1948-06-27 TMAX 86.0 NULL NULL G NULL
58	NZ000093844 1948-06-27 TMIN 4.0 NULL NULL G NULL
59	NZ000093844 1948-06-28 TMAX 122.0 NULL NULL G NULL
60	NZ000093844 1948-06-28 TMIN -8.0 NULL NULL G NULL
61	NZ000093844 1948-06-29 TMAX 89.0 NULL NULL G NULL
62	NZ000093844 1948-06-29 TMIN 19.0 NULL NULL G NULL
63	NZ000093844 1948-06-30 TMAX 72.0 NULL NULL G NULL
64	NZ000093844 1948-06-30 TMIN 22.0 NULL NULL G NULL
65	NZ000093844 1948-07-01 TMAX 73.0 NULL NULL G NULL
66	NZ000093844 1948-07-01 TMIN 43.0 NULL NULL G NULL
67	NZ000093844 1948-07-02 TMAX 88.0 NULL NULL G NULL
68	NZ000093844 1948-07-02 TMIN -6.0 NULL NULL G NULL
69	NZ000093844 1948-07-03 TMAX 82.0 NULL NULL G NULL
70	NZ000093844 1948-07-03 TMIN 7.0 NULL NULL G NULL
71	NZ000093844 1948-07-04 TMAX 106.0 NULL NULL G NULL
72	NZ000093844 1948-07-04 TMIN 17.0 NULL NULL G NULL
73	NZ000093844 1948-07-05 TMAX 121.0 NULL NULL G NULL
74	NZ000093844 1948-07-05 TMIN -11.0 NULL NULL G NULL
75	NZ000093844 1948-07-06 TMAX 133.0 NULL NULL G NULL
76	NZ000093844 1948-07-06 TMIN 23.0 NULL NULL G NULL
77	NZ000093844 1948-07-07 TMAX 106.0 NULL NULL G NULL
78	NZ000093844 1948-07-07 TMIN 42.0 NULL NULL G NULL
79	NZ000093844 1948-07-08 TMAX 106.0 NULL NULL G NULL
80	NZ000093844 1948-07-08 TMIN -22.0 NULL NULL G NULL
81	NZ000093844 1948-07-09 TMAX 112.0 NULL NULL G NULL
82	NZ000093844 1948-07-09 TMIN -13.0 NULL NULL G NULL
83	NZ000093844 1948-07-10 TMAX 77.0 NULL NULL G NULL
84	NZ000093844 1948-07-10 TMIN 21.0 NULL NULL G NULL
85	NZ000093844 1948-07-11 TMAX 97.0 NULL NULL G NULL
86	NZ000093844 1948-07-11 TMIN 36.0 NULL NULL G NULL
87	NZ000093844 1948-07-12 TMAX 129.0 NULL NULL G NULL
88	NZ000093844 1948-07-12 TMIN 11.0 NULL NULL G NULL
89	NZ000093844 1948-07-13 TMAX 106.0 NULL NULL G NULL
90	NZ000093844 1948-07-13 TMIN 37.0 NULL NULL G NULL
91	NZ000093844 1948-07-14 TMAX 84.0 NULL NULL G NULL
92	NZ000093844 1948-07-14 TMIN 36.0 NULL NULL G NULL
93	NZ000093844 1948-07-15 TMAX 96.0 NULL NULL G NULL
94	NZ000093844 1948-07-15 TMIN 4.0 NULL NULL G NULL
95	NZ000093844 1948-07-16 TMAX 107.0 NULL NULL G NULL
96	NZ000093844 1948-07-16 TMIN 8.0

104	NZ000093844	1948-07-20	TMIN	-19.0	NULL	NULL	G	NULL	
105	NZ000093844	1948-07-21	TMAX	78.0	NULL	NULL	G	NULL	
106	NZ000093844	1948-07-21	TMIN	18.0	NULL	NULL	G	NULL	
107	NZ000093844	1948-07-22	TMAX	97.0	NULL	NULL	G	NULL	
108	NZ000093844	1948-07-22	TMIN	31.0	NULL	NULL	G	NULL	
109	NZ000093844	1948-07-23	TMAX	89.0	NULL	NULL	G	NULL	
110	NZ000093844	1948-07-23	TMIN	64.0	NULL	NULL	G	NULL	
111	NZ000093844	1948-07-24	TMAX	101.0	NULL	NULL	G	NULL	
112	NZ000093844	1948-07-24	TMIN	52.0	NULL	NULL	G	NULL	
113	NZ000093844	1948-07-25	TMAX	111.0	NULL	NULL	G	NULL	
114	NZ000093844	1948-07-25	TMIN	9.0	NULL	NULL	G	NULL	
115	NZ000093844	1948-07-26	TMAX	122.0	NULL	NULL	G	NULL	
116	NZ000093844	1948-07-26	TMIN	-29.0	NULL	NULL	G	NULL	
117	NZ000093844	1948-07-27	TMAX	132.0	NULL	NULL	G	NULL	
118	NZ000093844	1948-07-27	TMIN	-34.0	NULL	NULL	G	NULL	
119	NZ000093844	1948-07-28	TMAX	106.0	NULL	NULL	G	NULL	
120	NZ000093844	1948-07-28	TMIN	-18.0	NULL	NULL	G	NULL	
121	NZ000093844	1948-07-29	TMAX	127.0	NULL	NULL	G	NULL	
122	NZ000093844	1948-07-29	TMIN	14.0	NULL	NULL	G	NULL	
123	NZ000093844	1948-07-30	TMAX	103.0	NULL	NULL	G	NULL	
124	NZ000093844	1948-07-30	TMIN	10.0	NULL	NULL	G	NULL	
125	NZ000093844	1948-07-31	TMAX	82.0	NULL	NULL	G	NULL	
126	NZ000093844	1948-07-31	TMIN	-1.0	NULL	NULL	G	NULL	
127	NZ000093844	1948-08-01	TMAX	77.0	NULL	NULL	G	NULL	
128	NZ000093844	1948-08-01	TMIN	-54.0	NULL	NULL	G	NULL	
129	NZ000093844	1948-08-02	TMAX	88.0	NULL	NULL	G	NULL	
130	NZ000093844	1948-08-02	TMIN	-36.0	NULL	NULL	G	NULL	
131	NZ000093844	1948-08-03	TMAX	131.0	NULL	NULL	G	NULL	
132	NZ000093844	1948-08-03	TMIN	34.0	NULL	NULL	G	NULL	
133	NZ000093844	1948-08-04	TMAX	120.0	NULL	NULL	G	NULL	
134	NZ000093844	1948-08-04	TMIN	17.0	NULL	NULL	G	NULL	
135	NZ000093844	1948-08-05	TMAX	134.0	NULL	NULL	G	NULL	
136	NZ000093844	1948-08-05	TMIN	36.0	NULL	NULL	G	NULL	
137	NZ000093844	1948-08-06	TMAX	116.0	NULL	NULL	G	NULL	
138	NZ000093844	1948-08-06	TMIN	23.0	NULL	NULL	G	NULL	
139	NZ000093844	1948-08-07	TMAX	116.0	NULL	NULL	G	NULL	
140	NZ000093844	1948-08-07	TMIN	31.0	NULL	NULL	G	NULL	
141	NZ000093844	1948-08-08	TMAX	127.0	NULL	NULL	G	NULL	
142	NZ000093844	1948-08-08	TMIN	39.0	NULL	NULL	G	NULL	
143	NZ000093844	1948-08-09	TMAX	146.0	NULL	NULL	G	NULL	
144	NZ000093844	1948-08-09	TMIN	23.0	NULL	NULL	G	NULL	
145	NZ000093844	1948-08-10	TMAX	145.0	NULL	NULL	G	NULL	
146	NZ000093844	1948-08-10	TMIN	23.0	NULL	NULL	G	NULL	
147	NZ000093844	1948-08-11	TMAX	116.0	NULL	NULL	G	NULL	
148	NZ000093844	1948-08-11	TMIN	-6.0	NULL	NULL	G	NULL	
149	NZ000093844	1948-08-12	TMAX	126.0	NULL	NULL	G	NULL	
150	NZ000093844	1948-08-12	TMIN	31.0	NULL	NULL	G	NULL	
151	NZ000093844	1948-08-13	TMAX	122.0	NULL	NULL	G	NULL	
152	NZ000093844	1948-08-13	TMIN	14.0	NULL	NULL	G	NULL	
153	NZ000093844	1948-08-14	TMAX	161.0	NULL	NULL	G	NULL	
154	NZ000093844	1948-08-14	TMIN	9.0	NULL	NULL	G	NULL	
155	NZ000093844	1948-08-15	TMAX	131.0	NULL	NULL	G	NULL	
156	NZ000093844	1948-08-15	TMIN	53.0	NULL	NULL	G	NULL	
157	NZ000093844	1948-08-16	TMAX	165.0	NULL	NULL	G	NULL	
158	NZ000093844	1948-08-16	TMIN	43.0	NULL	NULL	G	NULL	
159	NZ000093844	1948-08-17	TMAX	110.0	NULL	NULL	G	NULL	
160	NZ000093844	1948-08-17	TMIN	23.0	NULL	NULL	G	NULL	
161	NZ000093844	1948-08-18	TMAX	102.0	NULL	NULL	G	NULL	
162	NZ000093844	1948-08-18	TMIN	19.0	NULL	NULL	G	NULL	
163	NZ000093844	1948-08-19	TMAX	72.0	NULL	NULL	G	NULL	
164	NZ000093844	1948-08-19	TMIN	-37.0	NULL	NULL	G	NULL	
165	NZ000093844	1948-08-20	TMAX	84.0	NULL	NULL	G	NULL	
166	NZ000093844	1948-08-20	TMIN	-46.0	NULL	NULL	G	NULL	
167	NZ000093844	1948-08-21	TMAX	105.0	NULL	NULL	G	NULL	
168	NZ000093844	1948-08-21	TMIN	-57.0	NULL	NULL	G	NULL	
169	NZ000093844	1948-08-22	TMAX	119.0	NULL	NULL	G	NULL	
170	NZ000093844	1948-08-22	TMIN	-2.0	NULL	NULL	G	NULL	
171	NZ000093844	1948-08-23	TMAX	160.0	NULL	NULL	G	NULL	
172	NZ000093844	1948-08-23	TMIN	29.0	NULL	NULL	G	NULL	
173	NZ000093844	1948-08-24	TMAX	178.0	NULL	NULL	G	NULL	
174	NZ000093844	1948-08-24	TMIN	57.0	NULL	NULL	G	NULL	
175	NZ000093844	1948-08-25	TMAX	143.0	NULL	NULL	G	NULL	
176	NZ000093844	1948-08-25	TMIN	66.0	NULL	NULL	G	NULL	
177	NZ000093844	1948-08-26	TMAX	116.0	NULL	NULL	G	NULL	
178	NZ000093844	1948-08-26	TMIN	51.0	NULL	NULL	G	NULL	
179	NZ000093844	1948-08-27	TMAX	104.0	NULL	NULL	G	NULL	
180	NZ000093844	1948-08-27	TMIN	19.0	NULL	NULL	G	NULL	
181	NZ000093844	1948-08-28	TMAX	141.0	NULL	NULL	G	NULL	
182	NZ000093844	1948-08-28	TMIN	33.0	NULL	NULL	G	NULL	
183	NZ000093844	1948-08-29	TMAX	139.0	NULL	NULL	G	NULL	
184	NZ000093844	1948-08-29	TMIN	-19.0	NULL	NULL	G	NULL	
185	NZ000093844	1948-08-30	TMAX	129.0	NULL	NULL	G	NULL	
186	NZ000093844	1948-08-30	TMIN	8.0	NULL	NULL	G	NULL	



187	NZ000093844	1948-08-31	TMAX	128.0	NULL	NULL	G	NULL	
188	NZ000093844	1948-08-31	TMIN	-5.0	NULL	NULL	G	NULL	
189	NZ000093844	1948-09-01	TMAX	111.0	NULL	NULL	G	NULL	
190	NZ000093844	1948-09-01	TMIN	-31.0	NULL	NULL	G	NULL	
191	NZ000093844	1948-09-02	TMAX	131.0	NULL	NULL	G	NULL	
192	NZ000093844	1948-09-02	TMIN	-3.0	NULL	NULL	G	NULL	
193	NZ000093844	1948-09-03	TMAX	89.0	NULL	NULL	G	NULL	
194	NZ000093844	1948-09-03	TMIN	60.0	NULL	NULL	G	NULL	
195	NZ000093844	1948-09-04	TMAX	96.0	NULL	NULL	G	NULL	
196	NZ000093844	1948-09-04	TMIN	17.0	NULL	NULL	G	NULL	
197	NZ000093844	1948-09-05	TMAX	130.0	NULL	NULL	G	NULL	
198	NZ000093844	1948-09-05	TMIN	24.0	NULL	NULL	G	NULL	
199	NZ000093844	1948-09-06	TMAX	162.0	NULL	NULL	G	NULL	
200	NZ000093844	1948-09-06	TMIN	68.0	NULL	NULL	G	NULL	
201	NZ000093844	1948-09-07	TMAX	139.0	NULL	NULL	G	NULL	
202	NZ000093844	1948-09-07	TMIN	1.0	NULL	NULL	G	NULL	
203	NZ000093844	1948-09-08	TMAX	127.0	NULL	NULL	G	NULL	
204	NZ000093844	1948-09-08	TMIN	42.0	NULL	NULL	G	NULL	
205	NZ000093844	1948-09-09	TMAX	177.0	NULL	NULL	G	NULL	
206	NZ000093844	1948-09-09	TMIN	47.0	NULL	NULL	G	NULL	
207	NZ000093844	1948-09-10	TMAX	115.0	NULL	NULL	G	NULL	
208	NZ000093844	1948-09-10	TMIN	56.0	NULL	NULL	G	NULL	
209	NZ000093844	1948-09-11	TMAX	128.0	NULL	NULL	G	NULL	
210	NZ000093844	1948-09-11	TMIN	5.0	NULL	NULL	G	NULL	
211	NZ000093844	1948-09-12	TMAX	144.0	NULL	NULL	G	NULL	
212	NZ000093844	1948-09-12	TMIN	17.0	NULL	NULL	G	NULL	
213	NZ000093844	1948-09-13	TMAX	144.0	NULL	NULL	G	NULL	
214	NZ000093844	1948-09-13	TMIN	55.0	NULL	NULL	G	NULL	
215	NZ000093844	1948-09-14	TMAX	127.0	NULL	NULL	G	NULL	
216	NZ000093844	1948-09-14	TMIN	19.0	NULL	NULL	G	NULL	
217	NZ000093844	1948-09-15	TMAX	161.0	NULL	NULL	G	NULL	
218	NZ000093844	1948-09-15	TMIN	-1.0	NULL	NULL	G	NULL	
219	NZ000093844	1948-09-16	TMAX	149.0	NULL	NULL	G	NULL	
220	NZ000093844	1948-09-16	TMIN	41.0	NULL	NULL	G	NULL	
221	NZ000093844	1948-09-17	TMAX	152.0	NULL	NULL	G	NULL	
222	NZ000093844	1948-09-17	TMIN	58.0	NULL	NULL	G	NULL	
223	NZ000093844	1948-09-18	TMAX	149.0	NULL	NULL	G	NULL	
224	NZ000093844	1948-09-18	TMIN	39.0	NULL	NULL	G	NULL	
225	NZ000093844	1948-09-19	TMAX	136.0	NULL	NULL	G	NULL	
226	NZ000093844	1948-09-19	TMIN	39.0	NULL	NULL	G	NULL	
227	NZ000093844	1948-09-20	TMAX	98.0	NULL	NULL	G	NULL	
228	NZ000093844	1948-09-20	TMIN	-16.0	NULL	NULL	G	NULL	
229	NZ000093844	1948-09-21	TMAX	121.0	NULL	NULL	G	NULL	
230	NZ000093844	1948-09-21	TMIN	-12.0	NULL	NULL	G	NULL	
231	NZ000093844	1948-09-22	TMAX	186.0	NULL	NULL	G	NULL	
232	NZ000093844	1948-09-22	TMIN	95.0	NULL	NULL	G	NULL	
233	NZ000093844	1948-09-23	TMAX	156.0	NULL	NULL	G	NULL	
234	NZ000093844	1948-09-23	TMIN	48.0	NULL	NULL	G	NULL	
235	NZ000093844	1948-09-24	TMAX	124.0	NULL	NULL	G	NULL	
236	NZ000093844	1948-09-24	TMIN	9.0	NULL	NULL	G	NULL	
237	NZ000093844	1948-09-25	TMAX	161.0	NULL	NULL	G	NULL	
238	NZ000093844	1948-09-25	TMIN	44.0	NULL	NULL	G	NULL	
239	NZ000093844	1948-09-26	TMAX	130.0	NULL	NULL	G	NULL	
240	NZ000093844	1948-09-26	TMIN	57.0	NULL	NULL	G	NULL	
241	NZ000093844	1948-09-27	TMAX	144.0	NULL	NULL	G	NULL	
242	NZ000093844	1948-09-27	TMIN	61.0	NULL	NULL	G	NULL	
243	NZ000093844	1948-09-28	TMAX	119.0	NULL	NULL	G	NULL	
244	NZ000093844	1948-09-28	TMIN	21.0	NULL	NULL	G	NULL	
245	NZ000093844	1948-09-29	TMAX	94.0	NULL	NULL	G	NULL	
246	NZ000093844	1948-09-29	TMIN	17.0	NULL	NULL	G	NULL	
247	NZ000093844	1948-09-30	TMAX	82.0	NULL	NULL	G	NULL	
248	NZ000093844	1948-09-30	TMIN	27.0	NULL	NULL	G	NULL	
249	NZ000093844	1948-10-01	TMAX	83.0	NULL	NULL	G	NULL	
250	NZ000093844	1948-10-01	TMIN	31.0	NULL	NULL	G	NULL	
251	NZ000093844	1948-10-02	TMAX	104.0	NULL	NULL	G	NULL	
252	NZ000093844	1948-10-02	TMIN	34.0	NULL	NULL	G	NULL	
253	NZ000093844	1948-10-03	TMAX	99.0	NULL	NULL	G	NULL	
254	NZ000093844	1948-10-03	TMIN	-6.0	NULL	NULL	G	NULL	
255	NZ000093844	1948-10-04	TMAX	167.0	NULL	NULL	G	NULL	
256	NZ000093844	1948-10-04	TMIN	8.0	NULL	NULL	G	NULL	
257	NZ000093844	1948-10-05	TMAX	144.0	NULL	NULL	G	NULL	
258	NZ000093844	1948-10-05	TMIN	-27.0	NULL	NULL	G	NULL	
259	NZ000093844	1948-10-06	TMAX	114.0	NULL	NULL	G	NULL	
260	NZ000093844	1948-10-06	TMIN	47.0	NULL	NULL	G	NULL	
261	NZ000093844	1948-10-07	TMAX	128.0	NULL	NULL	G	NULL	
262	NZ000093844	1948-10-07	TMIN	57.0	NULL	NULL	G	NULL	
263	NZ000093844	1948-10-08	TMAX	115.0	NULL	NULL	G	NULL	
264	NZ000093844	1948-10-08	TMIN	36.0	NULL	NULL	G	NULL	
265	NZ000093844	1948-10-09	TMAX	124.0	NULL	NULL	G	NULL	
266	NZ000093844	1948-10-09	TMIN	61.0	NULL	NULL	G	NULL	
267	NZ000093844	1948-10-10	TMAX	111.0	NULL	NULL	G	NULL	
268	NZ000093844	1948-10-10	TMIN	-22.0	NULL	NULL	G	NULL	
269	NZ000093844	1948-10-11	TMAX	139.0	NULL	NULL	G	NULL	

270	NZ000093844 1948-10-11 TMIN 12.0 NULL NULL G NULL
271	NZ000093844 1948-10-12 TMAX 192.0 NULL NULL G NULL
272	NZ000093844 1948-10-12 TMIN 123.0 NULL NULL G NULL
273	NZ000093844 1948-10-13 TMAX 179.0 NULL NULL G NULL
274	NZ000093844 1948-10-13 TMIN 32.0 NULL NULL G NULL
275	NZ000093844 1948-10-14 TMAX 157.0 NULL NULL G NULL
276	NZ000093844 1948-10-14 TMIN 56.0 NULL NULL G NULL
277	NZ000093844 1948-10-15 TMAX 153.0 NULL NULL G NULL
278	NZ000093844 1948-10-15 TMIN 69.0 NULL NULL G NULL
279	NZ000093844 1948-10-16 TMAX 114.0 NULL NULL G NULL
280	NZ000093844 1948-10-16 TMIN 44.0 NULL NULL G NULL
281	NZ000093844 1948-10-17 TMAX 160.0 NULL NULL G NULL
282	NZ000093844 1948-10-17 TMIN 46.0 NULL NULL G NULL
283	NZ000093844 1948-10-18 TMAX 136.0 NULL NULL G NULL
284	NZ000093844 1948-10-18 TMIN 57.0 NULL NULL G NULL
285	NZ000093844 1948-10-19 TMAX 175.0 NULL NULL G NULL
286	NZ000093844 1948-10-19 TMIN 81.0 NULL NULL G NULL
287	NZ000093844 1948-10-20 TMAX 133.0 NULL NULL G NULL
288	NZ000093844 1948-10-20 TMIN 41.0 NULL NULL G NULL
289	NZ000093844 1948-10-21 TMAX 113.0 NULL NULL G NULL
290	NZ000093844 1948-10-21 TMIN 30.0 NULL NULL G NULL
291	NZ000093844 1948-10-22 TMAX 144.0 NULL NULL G NULL
292	NZ000093844 1948-10-22 TMIN 54.0 NULL NULL G NULL
293	NZ000093844 1948-10-23 TMAX 130.0 NULL NULL G NULL
294	NZ000093844 1948-10-23 TMIN 27.0 NULL NULL G NULL
295	NZ000093844 1948-10-24 TMAX 130.0 NULL NULL G NULL
296	NZ000093844 1948-10-24 TMIN 33.0 NULL NULL G NULL
297	NZ000093844 1948-10-25 TMAX 97.0 NULL NULL G NULL
298	NZ000093844 1948-10-25 TMIN 32.0 NULL NULL G NULL
299	NZ000093844 1948-10-26 TMAX 107.0 NULL NULL G NULL
300	NZ000093844 1948-10-26 TMIN 74.0 NULL NULL G NULL
301	NZ000093844 1948-10-27 TMAX 154.0 NULL NULL G NULL
302	NZ000093844 1948-10-27 TMIN 82.0 NULL NULL G NULL
303	NZ000093844 1948-10-28 TMAX 223.0 NULL NULL G NULL
304	NZ000093844 1948-10-28 TMIN 98.0 NULL NULL G NULL
305	NZ000093844 1948-10-29 TMAX 234.0 NULL NULL G NULL
306	NZ000093844 1948-10-29 TMIN 95.0 NULL NULL G NULL
307	NZ000093844 1948-10-30 TMAX 158.0 NULL NULL G NULL
308	NZ000093844 1948-10-30 TMIN 87.0 NULL NULL G NULL
309	NZ000093844 1948-10-31 TMAX 177.0 NULL NULL G NULL
310	NZ000093844 1948-10-31 TMIN 104.0 NULL NULL G NULL
311	NZ000093844 1948-11-01 TMAX 179.0 NULL NULL G NULL
312	NZ000093844 1948-11-01 TMIN 88.0 NULL NULL G NULL
313	NZ000093844 1948-11-02 TMAX 117.0 NULL NULL G NULL
314	NZ000093844 1948-11-02 TMIN -6.0 NULL NULL G NULL
315	NZ000093844 1948-11-03 TMAX 143.0 NULL NULL G NULL
316	NZ000093844 1948-11-03 TMIN 38.0 NULL NULL G NULL
317	NZ000093844 1948-11-04 TMAX 175.0 NULL NULL G NULL
318	NZ000093844 1948-11-04 TMIN 86.0 NULL NULL G NULL
319	NZ000093844 1948-11-05 TMAX 135.0 NULL NULL G NULL
320	NZ000093844 1948-11-05 TMIN 42.0 NULL NULL G NULL
321	NZ000093844 1948-11-06 TMAX 139.0 NULL NULL G NULL
322	NZ000093844 1948-11-06 TMIN 63.0 NULL NULL G NULL
323	NZ000093844 1948-11-07 TMAX 199.0 NULL NULL G NULL
324	NZ000093844 1948-11-07 TMIN 111.0 NULL NULL G NULL
325	NZ000093844 1948-11-08 TMAX 238.0 NULL NULL G NULL
326	NZ000093844 1948-11-08 TMIN 124.0 NULL NULL G NULL
327	NZ000093844 1948-11-09 TMAX 161.0 NULL NULL G NULL
328	NZ000093844 1948-11-09 TMIN 94.0 NULL NULL G NULL
329	NZ000093844 1948-11-10 TMAX 144.0 NULL NULL G NULL
330	NZ000093844 1948-11-10 TMIN 71.0 NULL NULL G NULL
331	NZ000093844 1948-11-11 TMAX 121.0 NULL NULL G NULL
332	NZ000093844 1948-11-11 TMIN 20.0 NULL NULL G NULL
333	NZ000093844 1948-11-12 TMAX 131.0 NULL NULL G NULL
334	NZ000093844 1948-11-12 TMIN 32.0 NULL NULL G NULL
335	NZ000093844 1948-11-13 TMAX 128.0 NULL NULL G NULL
336	NZ000093844 1948-11-13 TMIN 46.0 NULL NULL G NULL
337	NZ000093844 1948-11-14 TMAX 122.0 NULL NULL G NULL
338	NZ000093844 1948-11-14 TMIN 19.0 NULL NULL G NULL
339	NZ000093844 1948-11-15 TMAX 138.0 NULL NULL G NULL
340	NZ000093844 1948-11-15 TMIN 39.0 NULL NULL G NULL
341	NZ000093844 1948-11-16 TMAX 111.0 NULL NULL G NULL
342	NZ000093844 1948-11-16 TMIN 46.0 NULL NULL G NULL
343	NZ000093844 1948-11-17 TMAX 133.

353	NZ000093844 1948-11-22 TMAX	142.0 NULL	NULL	G	NULL	
354	NZ000093844 1948-11-22 TMIN	29.0 NULL	NULL	G	NULL	
355	NZ000093844 1948-11-23 TMAX	154.0 NULL	NULL	G	NULL	
356	NZ000093844 1948-11-23 TMIN	22.0 NULL	NULL	G	NULL	
357	NZ000093844 1948-11-24 TMAX	134.0 NULL	NULL	G	NULL	
358	NZ000093844 1948-11-24 TMIN	-6.0 NULL	NULL	G	NULL	
359	NZ000093844 1948-11-25 TMAX	180.0 NULL	NULL	G	NULL	
360	NZ000093844 1948-11-25 TMIN	6.0 NULL	NULL	G	NULL	
361	NZ000093844 1948-11-26 TMAX	228.0 NULL	NULL	G	NULL	
362	NZ000093844 1948-11-26 TMIN	79.0 NULL	NULL	G	NULL	
363	NZ000093844 1948-11-27 TMAX	150.0 NULL	NULL	G	NULL	
364	NZ000093844 1948-11-27 TMIN	0.0 NULL	NULL	G	NULL	
365	NZ000093844 1948-11-28 TMAX	203.0 NULL	NULL	G	NULL	
366	NZ000093844 1948-11-28 TMIN	1.0 NULL	NULL	G	NULL	
367	NZ000093844 1948-11-29 TMAX	201.0 NULL	NULL	G	NULL	
368	NZ000093844 1948-11-29 TMIN	64.0 NULL	NULL	G	NULL	
369	NZ000093844 1948-11-30 TMAX	134.0 NULL	NULL	G	NULL	
370	NZ000093844 1948-11-30 TMIN	39.0 NULL	NULL	G	NULL	
371	NZ000093844 1948-12-01 TMAX	117.0 NULL	NULL	G	NULL	
372	NZ000093844 1948-12-01 TMIN	44.0 NULL	NULL	G	NULL	
373	NZ000093844 1948-12-02 TMAX	128.0 NULL	NULL	G	NULL	
374	NZ000093844 1948-12-02 TMIN	52.0 NULL	NULL	G	NULL	
375	NZ000093844 1948-12-03 TMAX	147.0 NULL	NULL	G	NULL	
376	NZ000093844 1948-12-03 TMIN	92.0 NULL	NULL	G	NULL	
377	NZ000093844 1948-12-04 TMAX	142.0 NULL	NULL	G	NULL	
378	NZ000093844 1948-12-04 TMIN	49.0 NULL	NULL	G	NULL	
379	NZ000093844 1948-12-05 TMAX	172.0 NULL	NULL	G	NULL	
380	NZ000093844 1948-12-05 TMIN	82.0 NULL	NULL	G	NULL	
381	NZ000093844 1948-12-06 TMAX	223.0 NULL	NULL	G	NULL	
382	NZ000093844 1948-12-06 TMIN	95.0 NULL	NULL	G	NULL	
383	NZ000093844 1948-12-07 TMAX	177.0 NULL	NULL	G	NULL	
384	NZ000093844 1948-12-07 TMIN	82.0 NULL	NULL	G	NULL	
385	NZ000093844 1948-12-08 TMAX	184.0 NULL	NULL	G	NULL	
386	NZ000093844 1948-12-08 TMIN	63.0 NULL	NULL	G	NULL	
387	NZ000093844 1948-12-09 TMAX	150.0 NULL	NULL	G	NULL	
388	NZ000093844 1948-12-09 TMIN	7.0 NULL	NULL	G	NULL	
389	NZ000093844 1948-12-10 TMAX	198.0 NULL	NULL	G	NULL	
390	NZ000093844 1948-12-10 TMIN	53.0 NULL	NULL	G	NULL	
391	NZ000093844 1948-12-11 TMAX	243.0 NULL	NULL	G	NULL	
392	NZ000093844 1948-12-11 TMIN	115.0 NULL	NULL	G	NULL	
393	NZ000093844 1948-12-12 TMAX	198.0 NULL	NULL	G	NULL	
394	NZ000093844 1948-12-12 TMIN	144.0 NULL	NULL	G	NULL	
395	NZ000093844 1948-12-13 TMAX	189.0 NULL	NULL	G	NULL	
396	NZ000093844 1948-12-13 TMIN	77.0 NULL	NULL	G	NULL	
397	NZ000093844 1948-12-14 TMAX	162.0 NULL	NULL	G	NULL	
398	NZ000093844 1948-12-14 TMIN	49.0 NULL	NULL	G	NULL	
399	NZ000093844 1948-12-15 TMAX	128.0 NULL	NULL	G	NULL	
400	NZ000093844 1948-12-15 TMIN	56.0 NULL	NULL	G	NULL	
401	NZ000093844 1948-12-16 TMAX	138.0 NULL	NULL	G	NULL	
402	NZ000093844 1948-12-16 TMIN	78.0 NULL	NULL	G	NULL	
403	NZ000093844 1948-12-17 TMAX	172.0 NULL	NULL	G	NULL	
404	+-----+-----+-----+---+-----+-----+-----+-----+					
405	only showing top 400 rows					

## Data preparation for plot

```
1 pmonth_avg = (  
2     daily_nz_tmin_tmax.select("ID", "DATE", "ELEMENT", "VALUE")  
3     .withColumn("year", F.year("DATE"))  
4     .withColumn("month", F.month("DATE"))  
5     .groupBy("ID", "year", "month")  
6     .pivot("ELEMENT", ["TMIN", "TMAX"])  
7     .agg(F.avg("VALUE"))  
8     .withColumnRenamed("TMIN", "TMIN_avg")  
9     .withColumnRenamed("TMAX", "TMAX_avg")  
10 )  
11  
12 show_as_html(pmonth_avg)
```

```
1 /opt/spark/python/pyspark/sql/pandas/conversion.py:111: UserWarning: toPandas attempted Arrow optimization because 'spark.sql.execution.arrow.  
2 PyArrow >= 4.0.0 must be installed; however, it was not found.  
3 Attempting non-optimization as 'spark.sql.execution.arrow.pyspark.fallback.enabled' is set to true.  
4 warn(msg)
```

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

ID	year	month	TMIN_avg	TMAX_avg
NZM00093439	2006	10	91.933333	151.346154
NZ000093417	2006	11	106.500000	168.400000
NZ000093417	2006	9	94.210526	148.892857
NZM00093439	2014	10	79.941176	164.827586
NZM00093110	2014	2	152.500000	246.464286
NZ000093012	2003	7	66.838710	151.935484
NZ000093994	2003	2	206.214286	261.285714
NZM00093439	2003	3	129.736842	210.344828
NZ000933090	1985	10	83.774194	162.193548
NZ000937470	1985	8	-15.225806	94.645161
NZM00093781	1985	6	-2.333333	125.333333
NZM00093781	1985	2	108.680000	238.333333
NZM00093781	1985	4	57.947368	194.500000
NZ000093292	1993	7	42.900000	133.366667
NZ000937470	1993	10	55.923077	185.538462
NZ000093994	1971	3	199.322581	238.709677
NZM00093781	1965	8	20.076923	72.000000
NZ000939450	1950	11	44.233333	95.333333
NZ000933090	1950	3	99.677419	202.419355
NZ000093844	1950	2	77.607143	179.178571

## Plot monthly average TMIN, TMAX

```
1 | # ===== 0) Spark-Pandas: =====
2 | # first time run should install PyArrow to speed up .toPandas processing
3 | # !pip install PyArrow
4 | # spark.conf.set("spark.sql.execution.arrow.pyspark.enabled", "true")
5 | pdf = pmonth_avg.toPandas() # 列: ID, year, month, TMIN.avg, TMAX.avg
```

```
1 | /opt/spark/python/pyspark/sql/pandas/conversion.py:111: UserWarning: toPandas attempted Arrow optimization because 'spark.sql.execution.arrow.'
2 | PyArrow >= 4.0.0 must be installed; however, it was not found.
3 | Attempting non-optimization as 'spark.sql.execution.arrow.pyspark.fallback.enabled' is set to true.
4 | warn(msg)
```

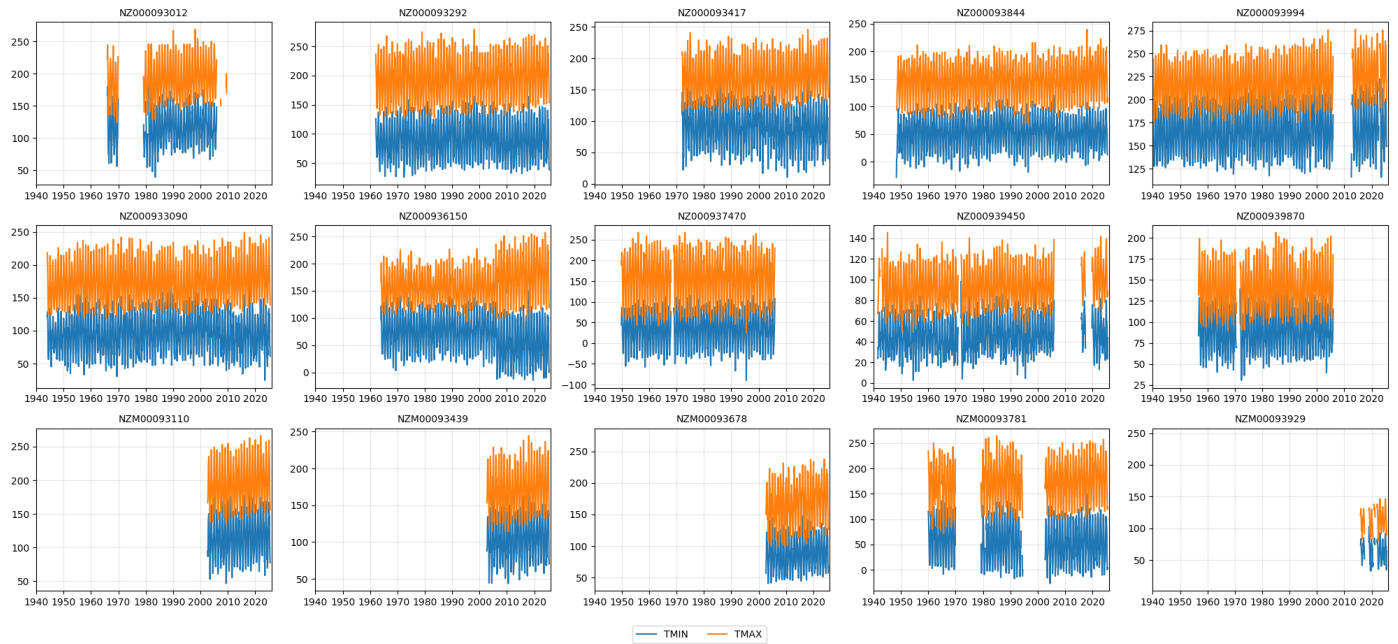
```
1 | # get the first day in each month
2 | pdf["month_start"] = pd.to_datetime(
3 |     pdf["year"].astype(str) + "-" + pdf["month"].astype(str) + "-01",
4 |     errors="coerce", # if conversion fails, the time will be set to Nat(Not a time).
5 | )
6 |
7 | # Keep only the columns needed for plotting & optimize types
8 | pdf = pdf[["ID", "month_start", "TMIN_avg", "TMAX_avg"]]
9 | pdf["ID"] = pdf["ID"].astype("category")
10 |
11 | # ===== 1) Alignment: reindex each station to a full monthly sequence (1940-01 ~ 2025-12) =====
12 | # Missing months are kept as NaN so that gaps remain visible in the plots
13 | full_index = pd.date_range("1940-01-01", "2025-12-01", freq="MS") # Month Start
14 |
15 |
16 | def reindex_station(g):
17 |     g = g.set_index("month_start")[["TMIN_avg", "TMAX_avg"]].reindex(
18 |         full_index
```

```

19     ) # Fill missing months with NaN
20     g.index.name = "month_start"
21     return g
22
23
24 blocks = []
25 for sid, g in pdf.groupby("ID", observed=True):
26     gg = reindex_station(g)
27     gg["ID"] = sid
28     blocks.append(gg.reset_index())
29
30 aligned = pd.concat(
31     blocks, ignore_index=True
32 ) # Final columns: month_start, TMIN_avg, TMAX_avg, ID
33 aligned = aligned[["ID", "month_start", "TMIN_avg", "TMAX_avg"]]
34
35 # ===== 2) Select 15 stations to plot (default: first 15, or replace with custom IDs) =====
36 ids_all = aligned["ID"].astype("category").cat.categories.tolist()
37 ids_15 = ids_all[:15] # Example customization: ids_15 = ["NZ000093012", "...", ...]
38
39 # ===== 3) Create 15 subplots (each station in its own axis, plot TMIN/TMAX with same x-range) =====
40 import matplotlib.pyplot as plt
41
42 nrows, ncols = 3, 5
43 fig, axes = plt.subplots(nrows=nrows, ncols=ncols, figsize=(20, 10), sharex=False)
44 axes = axes.ravel()
45
46 x_min = pd.Timestamp("1940-01-01")
47 x_max = pd.Timestamp("2025-12-31")
48
49 for ax, sid in zip(axes, ids_15):
50     g = aligned[aligned["ID"] == sid]
51     ax.plot(g["month_start"], g["TMIN_avg"], label="TMIN")
52     ax.plot(g["month_start"], g["TMAX_avg"], label="TMAX")
53     ax.set_title(str(sid), fontsize=10)
54     ax.set_xlim(x_min, x_max)
55     ax.grid(True, alpha=0.3)
56
57 # Hide unused subplots (in case there are fewer than 15 stations)
58 for k in range(len(ids_15), len(axes)):
59     axes[k].axis("off")
60
61 # Add a shared legend and overall layout adjustments
62 handles, labels = axes[0].get_legend_handles_labels()
63 fig.legend(handles, labels, loc="lower center", ncol=2)
64
65 fig.suptitle(
66     "NZ Stations • Monthly Mean TMIN/TMAX (NaN gaps, aligned 1940–2025)",
67     y=0.98,
68     fontsize=12,
69 )
70 fig.tight_layout(rect=[0, 0.04, 1, 0.96])
71 plt.show()
72
73 # Save figure
74 fig.savefig("./supplementary/nz_15stations_tmin_tmax_monthly_avg.png", dpi=300)

```

NZ Stations • Monthly Mean TMIN/TMAX (NaN gaps, aligned 1940–2025)



## Plot yearly average

```

1 pyear_avg = (
2     daily_nz_tmin_tmax.select("ID", "DATE", "ELEMENT", "VALUE")
3     .withColumn("year", F.year("DATE"))
4     .groupBy("ID", "year")
5     .pivot("ELEMENT", ["TMIN", "TMAX"])
6     .agg(F.avg("VALUE"))
7     .withColumnRenamed("TMIN", "TMIN_avg")
8     .withColumnRenamed("TMAX", "TMAX_avg")
9 )
10
11 show_as_html(pyear_avg)

```

```

1 /opt/spark/python/pyspark/sql/pandas/conversion.py:111: UserWarning: toPandas attempted Arrow optimization because 'spark.sql.execution.arrow.'
2   PyArrow >= 4.0.0 must be installed; however, it was not found.
3 Attempting non-optimization as 'spark.sql.execution.arrow.pyspark.fallback.enabled' is set to true.
4   warn(msg)

```

```

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

```

ID	year	TMIN_avg	TMAX_avg
NZ000093994	2019	169.269058	226.533742
NZ000093292	1999	98.085165	197.857534
NZ000093844	1991	51.408219	138.531507
NZM00093929	1994	13.000000	NaN
NZ000939450	2023	51.741697	106.269710
NZ000933090	2024	96.240741	179.644068
NZM00093929	2020	60.590909	102.609929
NZ000093012	1993	107.586301	186.632877
NZ000939870	1970	92.683824	143.194139
NZ000936150	1972	72.989071	154.704918

NZ000939450	1960	46.617486	95.576503
NZM00093929	2021	78.785714	127.712500
NZ000093012	1980	97.402597	181.987768
NZ000093012	1969	117.449367	178.551020
NZM00093439	2024	114.193309	177.675958
NZ000933090	1971	100.079452	183.115068
NZ000937470	1989	35.748603	160.530556
NZ000093844	1973	50.293151	142.676712
NZ000933090	1981	102.865385	178.986301
NZ000936150	1964	71.142077	149.901639

```

1 pdfy = pyear_avg.toPandas()
2 pdfy["ID"] = pdfy["ID"].astype("category")
3
4 # Create a datetime column representing the first day of each year
5 pdfy["year_start"] = pd.to_datetime(
6     pdfy["year"].astype(str) + "-01-01", errors="coerce"
7 )
8
9 # Keep only the columns needed for plotting
10 pdfy = pdfy[["ID", "year_start", "TMIN_avg", "TMAX_avg"]]
11
12 # Full yearly index from 1940 to 2025 (one timestamp per year start)
13 full_years = pd.date_range("1940-01-01", "2025-01-01", freq="YS")
14
15
16 def reindex_station_year(g):
17     # Reindex to the full yearly range; missing years become NaN
18     g2 = g.set_index("year_start")[["TMIN_avg", "TMAX_avg"]].reindex(full_years)
19     g2.index.name = "year_start"
20     return g2
21
22
23 blocks = []
24 for sid, g in pdfy.groupby("ID", observed=True):
25     gg = reindex_station_year(g)
26     gg["ID"] = sid
27     blocks.append(gg.reset_index())
28
29 aligned_year = pd.concat(blocks, ignore_index=True)
30 aligned_year = aligned_year[["ID", "year_start", "TMIN_avg", "TMAX_avg"]]
31
32
33 import matplotlib.pyplot as plt
34
35 # Select 15 stations (customize this list if needed)
36 ids_all = aligned_year["ID"].drop_duplicates().tolist()
37 ids_15 = ids_all[:15]
38
39 nrows, ncols = 3, 5
40 fig, axes = plt.subplots(nrows=nrows, ncols=ncols, figsize=(20, 10), sharex=False)
41 axes = axes.ravel()
42
43 x_min, x_max = pd.Timestamp("1940-01-01"), pd.Timestamp("2025-12-31")
44 # Optional: set a fixed y-axis for easier comparison across stations
45 # y_min, y_max = -5, 40
46
47 for ax, sid in zip(axes, ids_15):
48     g = aligned_year[aligned_year["ID"] == sid]
49     ax.plot(g["year_start"], g["TMIN_avg"], label="TMIN")
50     ax.plot(g["year_start"], g["TMAX_avg"], label="TMAX")
51     ax.set_title(str(sid), fontsize=10)
52     ax.set_xlim(x_min, x_max)
53     # ax.set_ylim(y_min, y_max)
54     ax.grid(True, alpha=0.3)
55
56 # Hide unused subplots if fewer than 15 stations
57 for k in range(len(ids_15), len(axes)):
58     axes[k].axis("off")
59
60 # Add a shared legend and adjust layout
61 handles, labels = axes[0].get_legend_handles_labels()

```

```

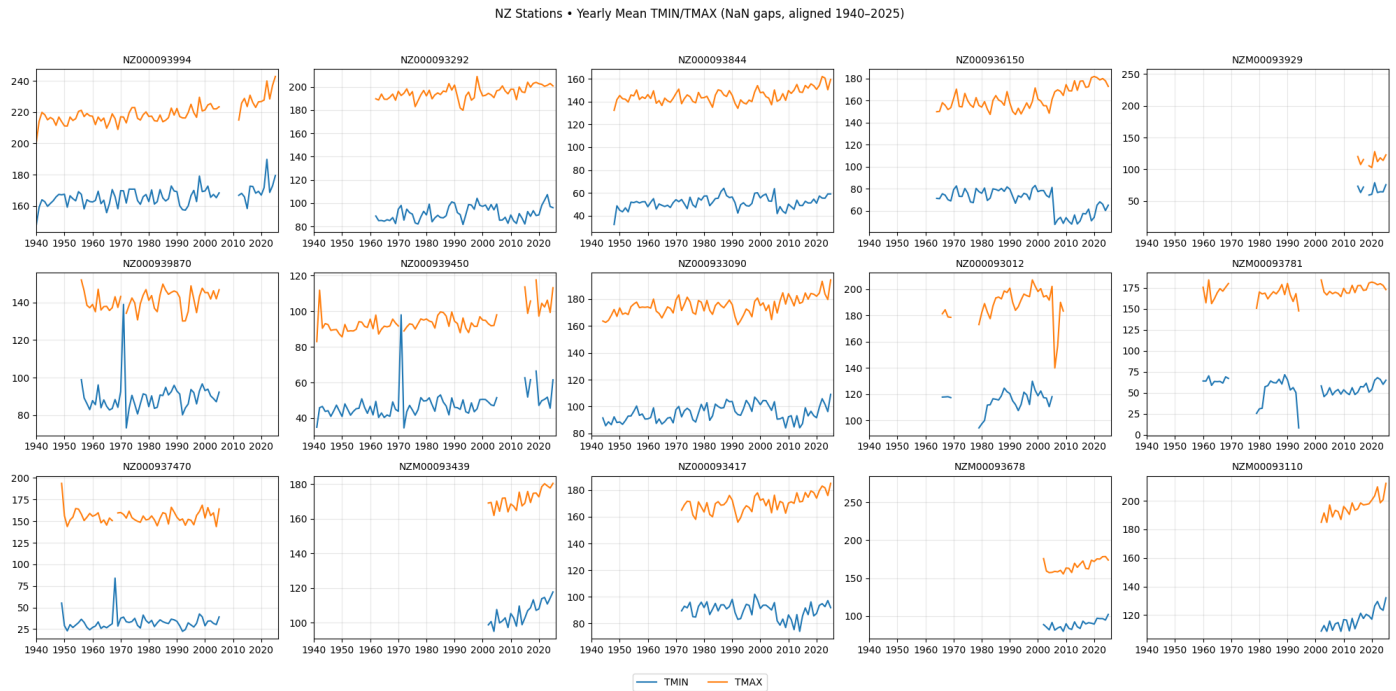
62 fig.legend(handles, labels, loc="lower center", ncol=2)
63 fig.suptitle(
64     "NZ Stations • Yearly Mean TMIN/TMAX (NaN gaps, aligned 1940–2025)",
65     y=0.98,
66     fontsize=12,
67 )
68 fig.tight_layout(rect=[0, 0.04, 1, 0.96])
69 plt.show()
70
71 # Save figure
72 fig.savefig("./supplementary/nz_15stations_tmin_tmax_yearly_avg.png", dpi=300)

```

```

1 /opt/spark/python/pyspark/sql/pandas/conversion.py:111: UserWarning: toPandas attempted Arrow optimization because 'spark.sql.execution.arrow.'
2 PyArrow >= 4.0.0 must be installed; however, it was not found.
3 Attempting non-optimization as 'spark.sql.execution.arrow.pyspark.fallback.enabled' is set to true.
4 warn(msg)

```



## Plot yearly average, filter out years that have less than 9 months data

```

1 # 0) Parameters
2 MIN_MONTHS = 9 # set to 10 if you want a stricter threshold
3 VALUE_IN_TENTHS = False # set True if raw VALUE is in tenths of °C
4
5
6 # 1) Spark — null invalid yearly means (variant B)
7
8 # 1) Monthly means per (ID, ELEMENT, year, month)
9 monthly_means = (
10     daily_nz_tmin_tmax.select("ID", "DATE", "ELEMENT", "VALUE")
11     .withColumn("DATE", F.to_date("DATE"))
12     .withColumn("year", F.year("DATE"))
13     .withColumn("month", F.month("DATE"))
14     .groupBy("ID", "ELEMENT", "year", "month")
15     .agg(F.avg("VALUE").alias("monthly_mean"))
16 )
17
18 # 2) Yearly stats per (ID, ELEMENT, year): count months and average of monthly means
19 year_stats = (
20     monthly_means.groupBy("ID", "ELEMENT", "year")
21     .agg(
22         F.count("*").alias("months_present"),
23         F.avg("monthly_mean").alias("year_mean_raw"),
24     )
25     .withColumn(
26         "year_mean",
27         F.when(
28             F.col("months_present") >= F.lit(MIN_MONTHS), F.col("year_mean_raw")
29         ).otherwise(F.lit(None)),
30     )

```



```

31 )
32
33 # 3) Pivot to wide format with NULLs where months_present < MIN_MONTHS
34 pyear_all = (
35     year_stats.groupBy("ID", "year")
36     .pivot("ELEMENT", ["TMIN", "TMAX"])
37     .agg(F.first("year_mean"))
38     .withColumnRenamed("TMIN", "TMIN_avg")
39     .withColumnRenamed("TMAX", "TMAX_avg")
40 )
41
42 # (Optional) diagnostics: months_present by element, for QC or annotation
43 months_diag = (
44     year_stats.groupBy("ID", "year")
45     .pivot("ELEMENT", ["TMIN", "TMAX"])
46     .agg(F.first("months_present"))
47     .withColumnRenamed("TMIN", "TMIN_months")
48     .withColumnRenamed("TMAX", "TMAX_months")
49 )
50
51 # Join diagnostics if you want them available downstream (safe to skip if not needed)
52 pyear_all = pyear_all.join(months_diag, on=["ID", "year"], how="left")
53
54
55 # 2) Pandas – align to full yearly index (1940–2025), keep NaN for gaps
56 import pandas as pd
57
58 spark.conf.set("spark.sql.execution.arrow.pyspark.enabled", "true")
59
60 pdfy = pyear_all.toPandas()
61
62 # Convert units if raw values are tenths of °C
63 if VALUE_IN_TENTHS:
64     for col in ["TMIN_avg", "TMAX_avg"]:
65         if col in pdfy.columns:
66             pdfy[col] = pdfy[col] / 10.0
67
68 pdfy["ID"] = pdfy["ID"].astype("category")
69 pdfy["year_start"] = pd.to_datetime(
70     pdfy["year"].astype(str) + "-01-01", errors="coerce"
71 )
72
73 # Keep only plotting columns (keep diagnostics if you want to annotate)
74 keep_cols = ["ID", "year_start", "TMIN_avg", "TMAX_avg"]
75 diag_cols = [c for c in ["TMIN_months", "TMAX_months"] if c in pdfy.columns]
76 pdfy = pdfy[keep_cols + diag_cols]
77
78 # Build full yearly index 1940–2025
79 full_years = pd.date_range("1940-01-01", "2025-01-01", freq="YS")
80
81
82 def reindex_station_year(g):
83     g2 = g.set_index("year_start").reindex(full_years)
84     g2.index.name = "year_start"
85     return g2
86
87
88 blocks = []
89 for sid, g in pdfy.groupby("ID", observed=True):
90     gg = reindex_station_year(g)
91     gg["ID"] = sid
92     blocks.append(gg.reset_index())
93
94 aligned_year = pd.concat(blocks, ignore_index=True)
95 aligned_year = aligned_year[["ID", "year_start", "TMIN_avg", "TMAX_avg"] + diag_cols]
96
97 # 3) Matplotlib – 15 subplots, TMIN/TMAX per station, aligned 1940–2025
98 import matplotlib.pyplot as plt
99
100 ids_all = aligned_year["ID"].drop_duplicates().tolist()
101 ids_15 = ids_all[:15] # customize this selection if needed
102
103 nrows, ncols = 3, 5
104 fig, axes = plt.subplots(nrows=nrows, ncols=ncols, figsize=(20, 10), sharex=False)
105 axes = axes.ravel()
106
107 x_min, x_max = pd.Timestamp("1940-01-01"), pd.Timestamp("2025-12-31")
108 # Optional fixed y-limits for comparability across stations (adjust ranges as needed)
109 # y_min, y_max = -5, 40
110
111 for ax, sid in zip(axes, ids_15):
112     g = aligned_year[aligned_year["ID"] == sid]
113     ax.plot(g["year_start"], g["TMIN_avg"], label="TMIN")

```

```

114     ax.plot(g["year_start"], g["TMAX_avg"], label="TMAX")
115     ax.set_title(str(sid), fontsize=10)
116     ax.set_xlim(x_min, x_max)
117     # ax.set_ylim(y_min, y_max)
118     ax.grid(True, alpha=0.3)
119
120 # Hide unused subplots if fewer than 15 stations
121 for k in range(len(ids_15), len(axes)):
122     axes[k].axis("off")
123
124 # Shared legend and layout
125 handles, labels = axes[0].get_legend_handles_labels()
126 fig.legend(handles, labels, loc="lower center", ncol=2)
127 fig.suptitle(
128     "NZ Stations • Yearly Mean TMIN/TMAX (NaN gaps via month coverage threshold, aligned 1940–2025)",
129     y=0.98,
130     fontsize=12,
131 )
132 fig.tight_layout(rect=[0, 0.04, 1, 0.96])
133
134 # Show and/or save
135 # plt.show()
136 fig.savefig(
137     "./supplementary/nz_15stations_tmin_tmax_yearly_avg_filtered_null.png", dpi=300
138 )

```

```

1 /opt/spark/python/pyspark/sql/pandas/conversion.py:111: UserWarning: toPandas attempted Arrow optimization because 'spark.sql.execution.arrow.'
2 PyArrow >= 4.0.0 must be installed; however, it was not found.
3 Attempting non-optimization as 'spark.sql.execution.arrow.pyspark.fallback.enabled' is set to true.
4 warn(msg)

```

NZ Stations • Yearly Mean TMIN/TMAX (NaN gaps via month coverage threshold, aligned 1940–2025)



```

1 nz_station_loc = stations_enriched.join(nz_station_ids, on="ID", how="inner")
2 show_as_html(nz_station_loc)

```

```

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
NZ000936150		NZ	-42.717	170.983	40.0	HOKITIKA AERODROME			93781

NZ000937470		NZ	-44.517	169.900	488.0	TARA HILLS	GSN		93747
NZ000939870		NZ	-43.950	-176.567	49.0	CHATHAM ISLANDS AWS			93987
NZM00093781		NZ	-43.489	172.532	37.5	CHRISTCHURCH INTL			93781
NZ000093012		NZ	-35.100	173.267	54.0	KAITAIA			93119
NZ000093292		NZ	-38.650	177.983	5.0	GISBORNE AERODROME	GSN		93292
NZM00093110		NZ	-37.000	174.800	7.0	AUCKLAND AERO AWS			93110
NZ000093417		NZ	-40.900	174.983	7.0	PARAPARAUMU AWS	GSN		93420
NZ000933090		NZ	-39.017	174.183	32.0	NEW PLYMOUTH AWS	GSN		93309
NZ000093994		NZ	-29.250	-177.917	49.0	RAOUL ISL/KERMADEC			93997
NZ000939450		NZ	-52.550	169.167	19.0	CAMPBELL ISLAND AWS	GSN		93947
NZM00093439		NZ	-41.333	174.800	12.0	WELLINGTON AERO AWS			93439
NZM00093929		NZ	-50.483	166.300	40.0	ENDERBY ISLAND AWS			93929
NZ000093844		NZ	-46.417	168.333	2.0	INVERCARGILL AIRPOR	GSN		93845
NZM00093678		NZ	-42.417	173.700	101.0	KAIKOURA			93678

15 rows × 21 columns

## NZ Station location plot

```
1  # === Imports ===
2  import geopandas as gpd
3  import matplotlib.pyplot as plt
4
5  # -----
6  # 1) Station table (Spark → pandas)
7  # -----
8  nz_pdf = nz_station_loc.select("ID", "LATITUDE", "LONGITUDE").toPandas()
9
10 # -----
11 # 2) Count distinct observation days per station
12 #   (remove duplicated TMIN/TMAX records on the same day)
13 # -----
14 obs_count_df = (
15     daily_nz_tmin_tmax.filter(F.col("ELEMENT").isin("TMIN", "TMAX"))
16     .select("ID", "DATE") # keep only station and date
17     .distinct() # remove duplicate TMIN/TMAX for the same day
18     .groupBy("ID")
19     .agg(F.count("*").alias("obs_count"))
20 )
21
22 obs_count_pdf = obs_count_df.toPandas()
23
24 # -----
25 # 3) Merge with station table & fill missing values
26 #   Ensure all 15 stations are included (fill NaN with 0)
27 # -----
28 merged = pd.merge(nz_pdf, obs_count_pdf, on="ID", how="left")
29 merged["obs_count"] = merged["obs_count"].fillna(0).astype(int)
30
31 # print(f"Stations listed: {len(nz_pdf)}")
32 # print(f"Stations after merge: {len(merged)}")
33 # if merged["obs_count"].isna().any():
34 #     print("Warning: some stations still have NaN in obs_count.")
35 # missing_ids = set(nz_pdf["ID"]) - set(merged["ID"])
36 # if missing_ids:
37 #     print("Missing station IDs in merged:", missing_ids)
38
39 # -----
40 # 4) Build GeoDataFrame (WGS84 → NZTM2000)
41 # -----
42 gdf_pts = gpd.GeoDataFrame(
43     merged,
44     geometry=gpd.points_from_xy(merged["LONGITUDE"], merged["LATITUDE"]),
45     crs="EPSG:4326",
46 )
47 gdf_pts_nztm = gdf_pts.to_crs(2193)
48
49 # Ensure all 15 stations are present
50 assert len(gdf_pts_nztm) == 15, f"Expected 15 stations, got {len(gdf_pts_nztm)}"
51
52 # -----
```

```

53 # 5) Load New Zealand polygon (Natural Earth) → NZTM2000
54 # -----
55 world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres"))
56 nz_poly = world[world["name"] == "New Zealand"].to_crs(2193)
57
58 # -----
59 # 6) Plotting (darker = more data → use 'Reds' colormap)
60 # -----
61 fig, ax = plt.subplots(figsize=(8, 10))
62
63 # 6a) Base map
64 nz_poly.plot(ax=ax, facecolor="#dfe8f2", edgecolor="#9fb5c8", linewidth=0.8, zorder=0)
65
66 # 6b) Station points (color by obs_count)
67 vals = gdf_pts_nztm["obs_count"].to_numpy(dtype=float)
68 vmin = float(np.nanmin(vals))
69 vmax = float(np.nanmax(vals)) if np.nanmax(vals) > 0 else 1.0
70
71 gdf_pts_nztm.plot(
72     ax=ax,
73     column="obs_count",
74     cmap="Reds", # darker = more data
75     markersize=80,
76     linewidth=0.4,
77     edgecolor="white",
78     legend=False,
79     vmin=vmin,
80     vmax=vmax,
81     zorder=1,
82 )
83
84 # 6c) Station labels
85 for _, r in gdf_pts_nztm.iterrows():
86     if r.geometry is not None:
87         ax.annotate(
88             r["ID"],
89             xy=(r.geometry.x, r.geometry.y),
90             xytext=(3, 3),
91             textcoords="offset points",
92             fontsize=6,
93             color="#404040",
94         )
95
96 # 6d) Map extent & title
97 minx, miny, maxx, maxy = nz_poly.total_bounds
98 pad_x = (maxx - minx) * 0.5
99 pad_y = (maxy - miny) * 0.5
100 ax.set_xlim(minx - pad_x, maxx + pad_x)
101 ax.set_ylim(miny - pad_y, maxy + pad_y)
102 ax.set_axis_off()
103
104 ax.set_title(
105     "Uneven Sample Sizes Across Stations (darker = more data): Implications for Nationwide Aggregation",
106     fontsize=12,
107 )
108
109 # 6e) Standalone colorbar (right side)
110 cax = fig.add_axes([0.92, 0.15, 0.02, 0.7])
111 sm = plt.cm.ScalarMappable(cmap="Reds", norm=plt.Normalize(vmin=vmin, vmax=vmax))
112 sm._A = []
113 cbar = fig.colorbar(sm, cax=cax)
114 cbar.set_label("Distinct Observation Days", fontsize=10)
115 plt.savefig(
116     "./supplementary/Uneven Sample Sizes Across Stations.png",
117     dpi=300,
118     bbox_inches="tight",
119 )
120 plt.show()

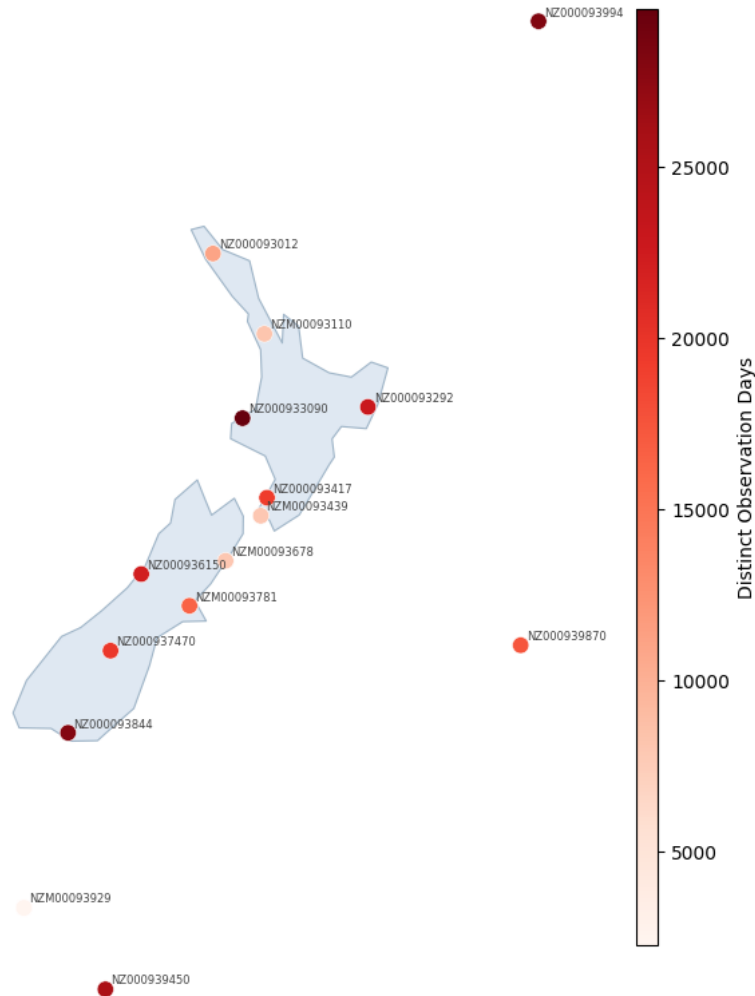
```

```

1 | /tmp/ipykernel_53/3388839370.py:55: FutureWarning: The geopandas.dataset module is deprecated and will be removed in GeoPandas 1.0. You can ge
2 | world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres"))

```

## Uneven Sample Sizes Across Stations (darker = more data): Implications for Nationwide Aggregation



## plot national region average TMIN TMAX

```
1 # -----
2 # NZ monthly area-weighted national means for TMIN / TMAX
3 # Time stratification (by month) + Voronoi area weights + availability re-normalization
4 # -----
5 #
6 # Algorithm overview & design considerations
7 # -----
8 # Goal:
9 #   Estimate New Zealand's national monthly means for TMIN and TMAX using
10 #   (1) time stratification (aggregate by month),
11 #   (2) spatial weights from a Voronoi tessellation over station locations, and
12 #   (3) per-month availability re-normalization so that months with missing stations
13 #       are not biased toward the subset that happened to report.
14 #
15 # Inputs:
16 #   - daily_nz_tmin_tmax (Spark DataFrame): columns ID, DATE (yyyy-mm-dd), ELEMENT in {TMIN,TMAX}, VALUE
17 #   - nz_station_loc      (Spark DataFrame): columns ID, LATITUDE, LONGITUDE (WGS84)
18 #
19 # Steps:
20 #   1) Pivot daily data to wide (TMIN/TMAX columns), then convert to pandas and add a month_start key.
21 #   2) For each station-month, compute daily-count QC and the monthly averages of TMIN/TMAX.
22 #       A station-month is considered valid if days_present >= MIN_DAYS_PER_MONTH.
23 #   3) Build spatial weights once, independent of time:
24 #       - Project station points to NZTM (EPSG:2193).
25 #       - Extract the New Zealand polygon from Natural Earth and project to EPSG:2193.
26 #       - Build a Voronoi tessellation with the NZ polygon as the envelope.
27 #       - Clip Voronoi cells to the NZ polygon.
28 #       - Map each cell to its nearest station; sum cell areas per station to get area weights.
29 #       - Normalize areas to sum to 1 to get w_i.
30 #   4) Merge the weights into the station-month table.
31 #   5) For each month, compute an availability re-normalized weighted mean:
32 #       - Filter to valid station-months having both a value and weight.
```

```

33 # - Re-normalize the weights within that month so the subset of available stations sums to 1.
34 # - Compute the weighted average for TMIN and TMAX separately.
35 # - Record how many stations contributed in that month (n_used_*).
36 #
37 # Why Voronoi weights?
38 # Voronoi area approximates each station's "zone of influence" given only point locations,
39 # offering a simple, reproducible spatial weighting when gridded climatologies are unavailable.
40 #
41 # CRS & geometry notes:
42 # - Areas are computed in EPSG:2193 (metres), so weights are physically meaningful.
43 # - We use Natural Earth's New Zealand polygon as the clipping boundary. For production,
44 #   replace with a higher-resolution NZ boundary if desired.
45 #
46 # QC & robustness:
47 # - MIN_DAYS_PER_MONTH guards against months with too few reports at a station.
48 # - The availability re-normalization removes bias when some stations are missing in a month.
49 # - If no valid stations exist in a month, the result is NaN.
50 #
51 # Caveats:
52 # - Voronoi-based areal weighting assumes spatial representativeness of stations; complex terrain,
53 #   microclimates, and coastal effects are not explicitly modeled.
54 # - Using Natural Earth boundaries may slightly mis-estimate coastal areas; swap in a better coastline if needed.
55 # - If stations move or IDs represent changing locations, re-compute weights per epoch or resolve metadata first.
56 #
57 # Output:
58 # - nz_monthly_avg (pandas DataFrame):
59 #   [month_start, TMIN_nzavg, TMAX_nzavg, n_used_TMIN, n_used_TMAX]
60 # - A PNG figure saved to ./supplementary/National_monthly_TM_avg.png
61 # -----
62
63 from shapely.ops import voronoi_diagram
64
65 # ===== 0) PARAMETERS =====
66 # Monthly validity threshold: a station-month is valid if it has at least this many daily observations
67 MIN_DAYS_PER_MONTH = 20
68
69 # ===== 1) LOAD / PREPARE TABULAR DATA =====
70 # daily_nz_tmin_tmax: Spark DataFrame with columns: ID, DATE (yyyy-mm-dd), ELEMENT in {TMIN,TMAX}, VALUE
71 # nz_station_loc: Spark DataFrame with columns: ID, LATITUDE, LONGITUDE
72
73 daily_nz_tmin_tmax = spark.read.parquet(daily_nz_tmin_tmax_path)
74
75 daily_nz_tm_wide = (
76     daily_nz_tmin_tmax.groupBy("ID", "DATE")
77     .pivot("ELEMENT", ["TMIN", "TMAX"])
78     .agg(F.first("VALUE"))
79 )
80
81 # Convert to pandas
82 daily_pdf = daily_nz_tm_wide.select("ID", "DATE", "TMIN", "TMAX").toPandas().copy()
83 stations_pdf = nz_station_loc.select("ID", "LATITUDE", "LONGITUDE").toPandas().copy()
84
85 # Parse dates & make month_start (month floor)
86 daily_pdf["DATE"] = pd.to_datetime(daily_pdf["DATE"])
87 daily_pdf["month_start"] = daily_pdf["DATE"].values.astype("datetime64[M]")
88
89 # ===== 2) STATION-MONTH AGGREGATION + QC =====
90 # Per-station per-month: mean TMIN/TMAX + number of distinct days present
91 monthly_station = daily_pdf.groupby(["ID", "month_start"], as_index=False).agg(
92     TMIN_avg=("TMIN", "mean"),
93     TMAX_avg=("TMAX", "mean"),
94     days_present=("DATE", "nunique"),
95 )
96 monthly_station["is_valid"] = monthly_station["days_present"] >= MIN_DAYS_PER_MONTH
97
98 # ===== 3) VORONOI AREA WEIGHTS (SPATIAL) =====
99 # 3.1 Stations -> GeoDataFrame (WGS84) and reproject to NZTM (EPSG:2193)
100 gdf_pts_wgs84 = gpd.GeoDataFrame(
101     stations_pdf,
102     geometry=gpd.points_from_xy(stations_pdf["LONGITUDE"], stations_pdf["LATITUDE"]),
103     crs="EPSG:4326",
104 )
105 gdf_pts_2193 = gdf_pts_wgs84.to_crs(2193)
106
107 # 3.2 New Zealand polygon (Natural Earth) -> EPSG:2193
108 world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres"))
109 nz_poly_2193 = world[world["name"] == "New Zealand"].to_crs(2193)
110 nz_envelope = nz_poly_2193.unary_union # used as Voronoi envelope and for clipping
111
112 # 3.3 Voronoi tessellation on the NZ envelope
113 # Note: shapely>=2.0 voronoi_diagram expects a MultiPoint / GeometryCollection
114 vor = voronoi_diagram(gdf_pts_2193.unary_union, envelope=nz_envelope)
115

```

```

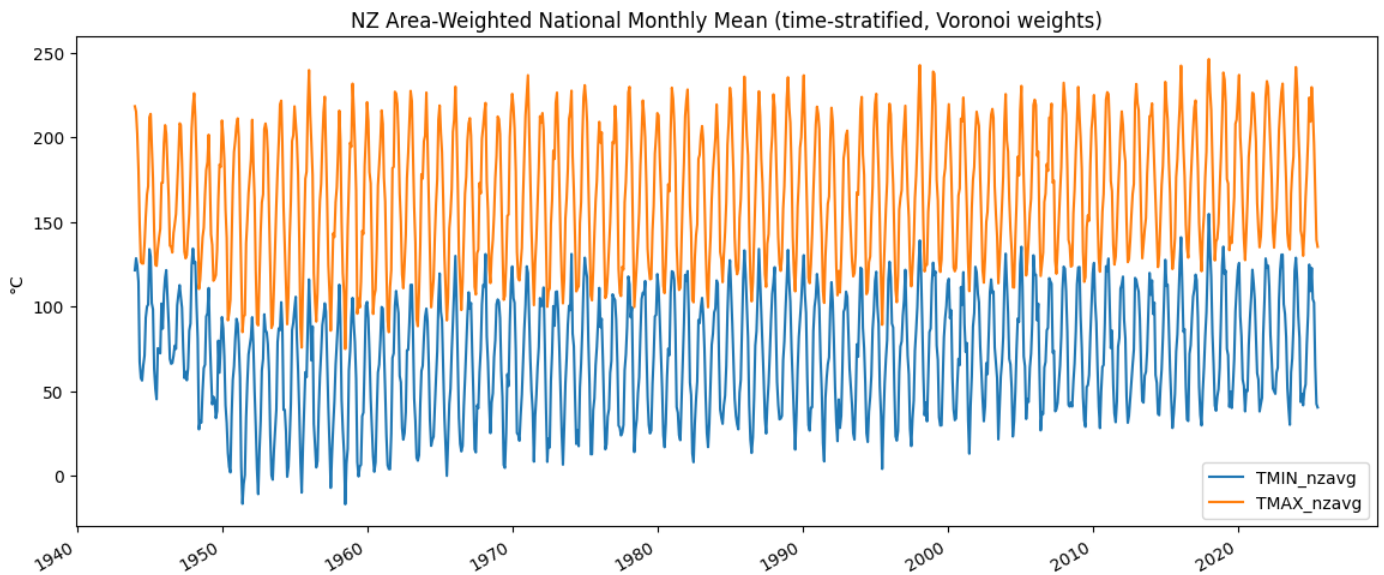
116 # Convert to GeoDataFrame and associate each cell with its nearest station
117 vor_gdf = gpd.GeoDataFrame(geometry=list(vor.geoms), crs=2193)
118
119 # Clip Voronoi cells to NZ boundary
120 vor_clip = gpd.overlay(vor_gdf, nz_poly_2193, how="intersection")
121
122 # Join each cell to the nearest station (each cell corresponds to one nearest station)
123 # Requires a spatial index (pygeos or rtree) to be installed for performance.
124 cell_to_station = vor_clip.sjoin_nearest(gdf_pts_2193[["ID", "geometry"]], how="left")
125 cell_to_station = cell_to_station.rename(columns={"ID": "ID"})
126
127 # 3.4 Compute area weights (aggregate: a station can own multiple clipped fragments)
128 cell_to_station["area"] = cell_to_station.geometry.area
129 weights_df = (
130     cell_to_station.groupby("ID", as_index=False)["area"]
131     .sum()
132     .rename(columns={"area": "area_total"})
133 )
134
135 # Normalize to get per-station weights w_i (sum to 1 over all stations)
136 weights_df["w"] = weights_df["area_total"] / weights_df["area_total"].sum()
137
138 # ===== 4) MERGE WEIGHTS WITH STATION-MONTH TABLE =====
139 aligned = monthly_station.merge(weights_df[["ID", "w"]], on="ID", how="left")
140
141
142 # ===== 5) TIME-STRATIFIED, AVAILABILITY RE-NORMALIZED MEAN =====
143 def weighted_mean_in_month(g, col):
144     """
145     Compute the re-normalized Voronoi-weighted mean for a single month.
146     g: rows for one month
147     col: 'TMIN_avg' or 'TMAX_avg'
148     Returns (value, n_used) where n_used is the number of contributing stations.
149     """
150     g_valid = g[g["is_valid"] & g[col].notna() & g["w"].notna()]
151     if g_valid.empty:
152         return np.nan, 0
153     w_norm = g_valid["w"] / g_valid["w"].sum()
154     val = np.sum(w_norm * g_valid[col].to_numpy())
155     return val, len(g_valid)
156
157
158 rows = []
159 for month, g in aligned.groupby("month_start"):
160     tmin_val, n_tmin = weighted_mean_in_month(g, "TMIN_avg")
161     tmax_val, n_tmax = weighted_mean_in_month(g, "TMAX_avg")
162     rows.append(
163         {
164             "month_start": month,
165             "TMIN_nzavg": tmin_val,
166             "TMAX_nzavg": tmax_val,
167             "n_used_TMIN": n_tmin,
168             "n_used_TMAX": n_tmax,
169         }
170     )
171
172 nz_monthly_avg = pd.DataFrame(rows).sort_values("month_start").reset_index(drop=True)
173
174 # ===== 6) OPTIONAL: PLOT =====
175 import matplotlib.pyplot as plt
176
177 fig, ax = plt.subplots(figsize=(14, 6))
178 ax.plot(nz_monthly_avg["month_start"], nz_monthly_avg["TMIN_nzavg"], label="TMIN_nzavg")
179 ax.plot(nz_monthly_avg["month_start"], nz_monthly_avg["TMAX_nzavg"], label="TMAX_nzavg")
180 ax.set_title(
181     "NZ Area-Weighted National Monthly Mean (time-stratified, Voronoi weights)"
182 )
183 ax.set_ylabel("°C")
184 ax.legend()
185 fig.autofmt_xdate()
186
187 # Save a copy (adjust path as needed)
188 plt.savefig("./supplementary/National_monthly_TM_avg.png", dpi=300, bbox_inches="tight")
189 plt.show()
190
191 # Final table:
192 #   nz_monthly_avg → columns:
193 #       [month_start, TMIN_nzavg, TMAX_nzavg, n_used_TMIN, n_used_TMAX]

```

```

1 | /tmp/ipykernel_53/224542698.py:108: FutureWarning: The geopandas.dataset module is deprecated and will be removed in GeoPandas 1.0. You can ge
2 | world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres"))

```

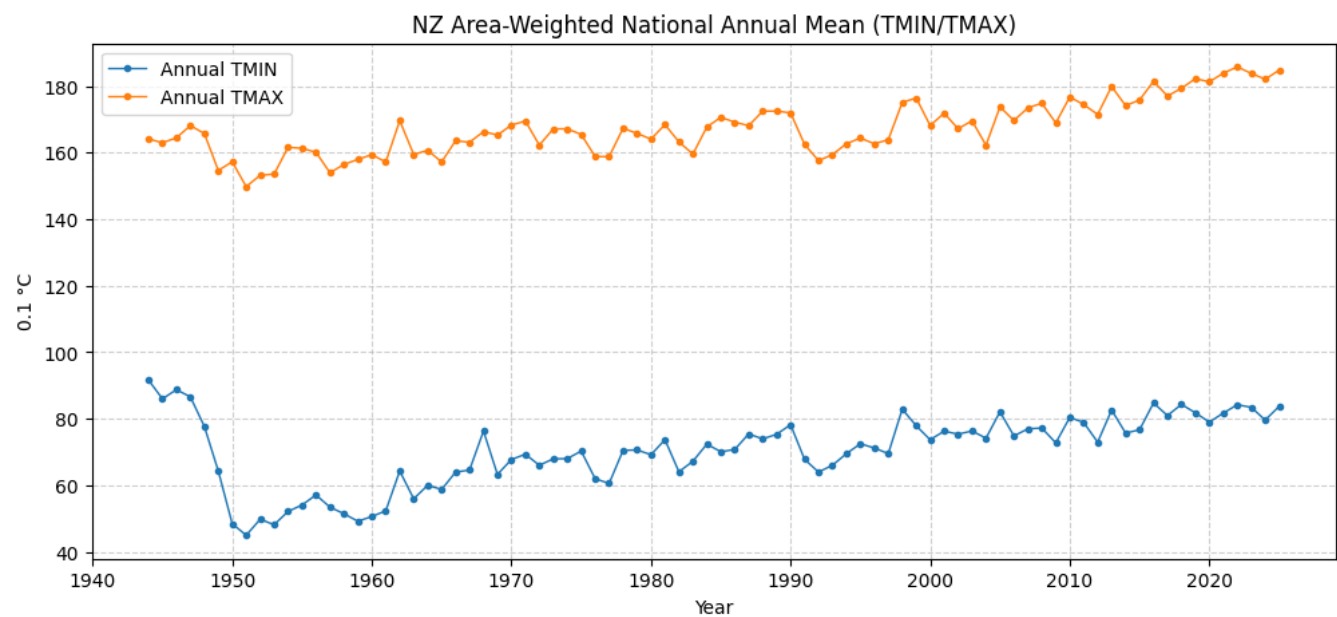


```

1 # ===== 7) AGGREGATE TO ANNUAL MEANS =====
2 # Add year column
3 nz_monthly_avg["year"] = nz_monthly_avg["month_start"].dt.year
4
5 # For each year, compute mean TMIN/TMAX across valid months
6 nz_annual_avg = nz_monthly_avg.groupby("year", as_index=False).agg(
7     TMIN_nzavg=("TMIN_nzavg", "mean"),
8     TMAX_nzavg=("TMAX_nzavg", "mean"),
9     n_used_TMIN=("n_used_TMIN", "sum"),
10    n_used_TMAX=("n_used_TMAX", "sum"),
11 )
12
13 # ===== 8) PLOT ANNUAL MEANS =====
14 fig, ax = plt.subplots(figsize=(12, 5))
15 ax.plot(
16     nz_annual_avg["year"],
17     nz_annual_avg["TMIN_nzavg"],
18     marker="o",
19     markersize=3,
20     linewidth=1,
21     label="Annual TMIN",
22 )
23 ax.plot(
24     nz_annual_avg["year"],
25     nz_annual_avg["TMAX_nzavg"],
26     marker="o",
27     markersize=3,
28     linewidth=1,
29     label="Annual TMAX",
30 )
31 ax.set_title("NZ Area-Weighted National Annual Mean (TMIN/TMAX)")
32 ax.set_ylabel("0.1 °C")
33 ax.set_xlabel("Year")
34 ax.legend()
35 plt.grid(True, linestyle="--", alpha=0.6)
36
37 plt.savefig("./supplementary/National_annual_TM_avg.png", dpi=300, bbox_inches="tight")
38 plt.show()
39
40 # Final table: nz_annual_avg → [year, TMIN_nzavg, TMAX_nzavg, n_used_TMIN, n_used_TMAX]

```





## Precipitation Plot

### PRCP station stats

```
1 # when processed, directly load from cloud.
2 daily_prdp = (
3     daily.filter(F.col("ELEMENT") == "PRCP")
4     .select(["ID", "DATE", "VALUE"])
5     .withColumn("DATE", F.to_date(F.col("DATE"), "yyyyMMdd"))
6     .withColumnRenamed("VALUE", "PRCP_VALUE")
7     .withColumn("year", F.year("DATE"))
8 )
9
10 prcp_stations = daily_prdp.join(
11     F.broadcast(stations_enriched.filter(F.col("PRCP") == "1")),
12     on="ID",
13     how="left",
14 ).withColumnRenamed("NAME", "STATION_NAME")
15
16 show_as_html(prcp_stations)
```

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

ID	DATE	PRCP_VALUE	year	STATE	COUNTRY_CODE	LATITUDE	LONGITUDE	ELEVATION	STATION_NAME
ASN00030018	2008-01-01	0.0	2008		AS	-18.2922	143.5483	291.7	GEORGETOWN POST OFFICE
ASN00030019	2008-01-01	0.0	2008		AS	-19.2647	143.6741	536.0	GILBERTON
ASN00030021	2008-01-01	0.0	2008		AS	-20.7389	144.4853	430.0	GLENDOWER STATION
ASN00030022	2008-01-01	0.0	2008		AS	-20.8192	144.2333	316.4	HUGHENDEN AIRPORT
ASN00030025	2008-01-01	0.0	2008		AS	-20.8528	144.2258	320.0	HUGHENDEN STATION
ASN00030009	2008-01-01	0.0	2008		AS	-20.0078	144.8992	722.4	CARGOON STATION
ASN00029121	2008-01-01	0.0	2008		AS	-20.5958	139.6939	280.0	WEST LEICHHARDT STATION
ASN00029123	2008-01-01	0.0	2008		AS	-20.5833	139.5767	330.0	LAKE MOONDARRA
ASN00029126	2008-01-01	0.0	2008		AS	-20.7361	139.4817	381.0	MOUNT ISA MINE
ASN00029127	2008-01-01	0.0	2008		AS	-20.6778	139.4875	340.3	MOUNT ISA AERO

ASN00029129	2008-01-01	0.0	2008		AS	-21.2150	140.2333	282.0	DEVONCOURT STATION
ASN00029130	2008-01-01	0.0	2008		AS	-20.9542	139.5861	339.0	MIM RIFLE CREEK
ASN00029131	2008-01-01	0.0	2008		AS	-20.0200	139.8433	190.0	GERETA STATION
ASN00029132	2008-01-01	0.0	2008		AS	-20.6556	142.1006	120.0	MANFRED DOWNS STATION
ASN00029136	2008-01-01	0.0	2008		AS	-21.3661	140.4986	300.0	FARLEY STATION
ASN00029137	2008-01-01	0.0	2008		AS	-19.6661	141.3961	100.0	NUMIL DOWNS STATION
ASN00029139	2008-01-01	180.0	2008		AS	-17.1142	139.5981	4.0	SWEERS ISLAND
ASN00029141	2008-01-01	0.0	2008		AS	-20.6664	140.5050	186.0	CLONCURRY AIRPORT
ASN00029144	2008-01-01	0.0	2008		AS	-20.4350	141.5481	120.0	KESWICK STATION
ASN00029150	2008-01-01	0.0	2008		AS	-20.8422	141.4292	150.0	MALVIE DOWNS STATION

20 rows × 24 columns

```
1 prcp_pdf = (  
2     prcp_stations.groupBy("year", "COUNTRY_CODE", "COUNTRY_NAME")  
3     .agg(F.avg("PRCP_VALUE").alias("PRCP_yavg"))  
4     .orderBy("PRCP_yavg")  
5 )  
6  
7  
8 show_as_html(prcp_pdf)  
9  
10 # when needed, make next line run.  
11 # prcp_pdf.write.mode("overwrite").parquet(prcp_pdf_path)
```

```
1 # when needed, make next line run.  
2 # prcp_pdf.write.mode("overwrite").parquet(prcp_pdf_path)  
3 !hdfs dfs -ls -h {prcp_pdf_path}
```

```
1 Found 2 items  
2 -rw-r--r--  1 yxi75 supergroup          0 2025-09-14 07:29 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/prcp_pdf_parquet/_S  
3 -rw-r--r--  1 yxi75 supergroup    194.3 K 2025-09-14 07:29 wasbs://campus-user@mdsstorage002.blob.core.windows.net/yxi75/prcp_pdf_parquet/pa
```

```
1 prcp_pdf = spark.read.parquet(prcp_pdf_path)  
2 prcp_pdf.cache()  
3 show_as_html(prcp_pdf)
```

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

year	COUNTRY_CODE	COUNTRY_NAME	PRCP_yavg
1874	UK	United Kingdom	-1.076712
2021	WQ	Wake Island [United States]	0.000000
2021	VE	Venezuela	0.000000
2021	AE	United Arab Emirates	0.000000
2021	BG	Bangladesh	0.000000
2002	MV	Maldives	0.000000
2002	CK	Cocos (Keeling) Islands [Australia]	0.000000
2020	WQ	Wake Island [United States]	0.000000
2004	BM	Burma	0.000000
2004	MV	Maldives	0.000000

2004	KT	Christmas Island [Australia]	0.000000
2004	LI	Liberia	0.000000
2001	SB	Saint Pierre and Miquelon [France]	0.000000
1990	MC	Macau S.A.R	0.000000
1990	BM	Burma	0.000000
1988	KU	Kuwait	0.000000
1988	QA	Qatar	0.000000
1988	MC	Macau S.A.R	0.000000
1978	CK	Cocos (Keeling) Islands [Australia]	0.000000
1978	AE	United Arab Emirates	0.000000

```
1 # Calculate summary statistics from the global precipitation average dataframe
2 prcp_stats = prcp_pdf.agg(
3     F.countDistinct("COUNTRY_CODE").alias("unique_country_codes"),
4     F.count("*").alias("total_rows"),
5     F.countDistinct("YEAR").alias("unique_years"),
6 ).collect()[0]
7
8 # Print the summary statistics with improved, more informative output strings
9 print(f"Number of unique countries: {prcp_stats['unique_country_codes']}")
10 print(f"Total number of records: {prcp_stats['total_rows']}")
11 print(f"Number of years in the dataset: {prcp_stats['unique_years']}")
```

```
1 [Stage 732:=====>(106 + 1) / 107]
2
3 There are 17731 PRCP records.
```

```
1 prcp_pdf = prcp_pdf.orderBy(F.asc("PRCP_yavg")).toPandas()
```

```
1 # Calculate descriptive statistics grouped by year and country
2 stats_by_year_country = prcp_pdf.groupby(["year", "COUNTRY_NAME"])[
3     "PRCP_yavg"
4 ].describe()
5
6 # Display grouped statistics
7 print("Descriptive statistics for average rainfall by year and country:")
8 print(stats_by_year_country)
9
10 # Calculate and display overall descriptive statistics
11 overall_stats = prcp_pdf["PRCP_yavg"].describe()
12 print("\nOverall average rainfall descriptive statistics:")
13 print(overall_stats)
```

1	Descriptive statistics for average rainfall by year and country:				
2		count	mean	std	min \
3	year COUNTRY_NAME				
4	1750 Australia	1.0	23.187021	NaN	23.187021
5	1781 Germany	1.0	24.558904	NaN	24.558904
6	1782 Germany	1.0	13.712329	NaN	13.712329
7	1783 Germany	1.0	17.832877	NaN	17.832877
8	1784 Germany	1.0	16.576503	NaN	16.576503
9	...	...	...	...	...
10	2025 Virgin Islands [United States]	1.0	27.814272	NaN	27.814272
11	Wake Island [United States]	1.0	21.909091	NaN	21.909091
12	Wallis and Futuna [France]	1.0	151.692308	NaN	151.692308
13	Zambia	1.0	59.500000	NaN	59.500000
14	Zimbabwe	1.0	54.647668	NaN	54.647668
15					
16		25%	50%	75%	\
17	year COUNTRY_NAME				
18	1750 Australia	23.187021	23.187021	23.187021	
19	1781 Germany	24.558904	24.558904	24.558904	
20	1782 Germany	13.712329	13.712329	13.712329	
21	1783 Germany	17.832877	17.832877	17.832877	
22	1784 Germany	16.576503	16.576503	16.576503	
23	...	...	...	...	
24	2025 Virgin Islands [United States]	27.814272	27.814272	27.814272	
25	Wake Island [United States]	21.909091	21.909091	21.909091	

```
26 | Wallis and Futuna [France] 151.692308 151.692308 151.692308
27 | Zambia 59.500000 59.500000 59.500000
28 | Zimbabwe 54.647668 54.647668 54.647668
29 |
30 | max
31 | year COUNTRY_NAME
32 | 1750 Australia 23.187021
33 | 1781 Germany 24.558904
34 | 1782 Germany 13.712329
35 | 1783 Germany 17.832877
36 | 1784 Germany 16.576503
37 | ... ...
38 | 2025 Virgin Islands [United States] 27.814272
39 | Wake Island [United States] 21.909091
40 | Wallis and Futuna [France] 151.692308
41 | Zambia 59.500000
42 | Zimbabwe 54.647668
43 |
44 | [17726 rows x 8 columns]
45 |
46 | Overall average rainfall descriptive statistics:
47 | count 17731.000000
48 | mean 44.393026
49 | std 198.491195
50 | min -1.076712
51 | 25% 15.557789
52 | 50% 25.204204
53 | 75% 48.160650
54 | max 15875.000000
55 | Name: PRCP_yavg, dtype: float64
```

```
1 | # the highest average rainfl in a single year country
2 | max_prdp = prcp_pdf.loc[prcp_pdf.PRCP_yavg == prcp_pdf.PRCP_yavg.max()]
3 | print(max_prdp)
```

```
1 | year COUNTRY_CODE COUNTRY_NAME PRCP_yavg
2 | 17730 1952 None None 15875.0
```

```
1 | # how many prcp recorded stations are there?
2 | # stations:127610 /129657 (prcp num/total num)
3 | prcp_station_ids = daily_prdp.select("ID").distinct()
4 | prcp_station_ids_count = prcp_station_ids.count()
5 | print(f"There are {prcp_station_ids_count} stations which recored PRCP value.")
```

```
1 | [Stage 13:=====>(106 + 1) / 107]
2 |
3 | There are 127610 stations which recored PRCP value.
```

```
1 | # how many countries are these prcp stations distributed?
2 | # countries: 218/219 ((prcp num/total num))
3 | prcp_station_country_count = prcp_stations.select("COUNTRY_NAME").distinct().count()
4 | print(f"These 127610 PRCP stations distribute in {prcp_station_country_count} countries.")
```

## outliers study

```
1 | show_as_html(prcp_pdf)
```

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

year	COUNTRY_CODE	COUNTRY_NAME	PRCP_yavg
1874	UK	United Kingdom	-1.076712
2021	WQ	Wake Island [United States]	0.000000

2021	VE	Venezuela	0.000000
2021	AE	United Arab Emirates	0.000000
2021	BG	Bangladesh	0.000000
2002	MV	Maldives	0.000000
2002	CK	Cocos (Keeling) Islands [Australia]	0.000000
2020	WQ	Wake Island [United States]	0.000000
2004	BM	Burma	0.000000
2004	MV	Maldives	0.000000
2004	KT	Christmas Island [Australia]	0.000000
2004	LI	Liberia	0.000000
2001	SB	Saint Pierre and Miquelon [France]	0.000000
1990	MC	Macau S.A.R	0.000000
1990	BM	Burma	0.000000
1988	KU	Kuwait	0.000000
1988	QA	Qatar	0.000000
1988	MC	Macau S.A.R	0.000000
1978	CK	Cocos (Keeling) Islands [Australia]	0.000000
1978	AE	United Arab Emirates	0.000000

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 import numpy as np
4
5 # Use the PRCP_yavg column from prcp_pdf
6 data = (
7     prcp_pdf
8     .select("PRCP_yavg")
9     .where(F.col("PRCP_yavg").isNotNull())
10    .toPandas()["PRCP_yavg"]
11 )
12
13 # Threshold for clipping (99th percentile)
14 clip_threshold = np.percentile(data, 99)
15
16 # Create figure with 2 subplots (stacked vertically)
17 fig, axes = plt.subplots(2, 1, figsize=(12, 8))
18
19 # --- (A) Boxplot with clipping (focus on bulk distribution) ---
20 sns.boxplot(
21     x=data[data < clip_threshold],
22     ax=axes[0],
23     color="skyblue",
24     fliersize=2
25 )
26 axes[0].set_title("Distribution of Annual Average Rainfall (<99th Percentile)")
27 axes[0].set_xlabel("Average Rainfall")
28
29 # --- (B) Boxplot with log scale (full data, outliers compressed) ---
30 sns.boxplot(
31     x=np.log1p(data), # log(1+x) transform
32     ax=axes[1],
33     color="lightgreen",
34     fliersize=2
35 )
36 axes[1].set_title("Distribution of Annual Average Rainfall (Log Scale)")
37 axes[1].set_xlabel("log(1 + Average Rainfall)")
38
39 plt.tight_layout()
40
41 plt.savefig("./supplementary/annual_PRCP_outlier_bboxplot.png", dpi=220)
42 plt.show()

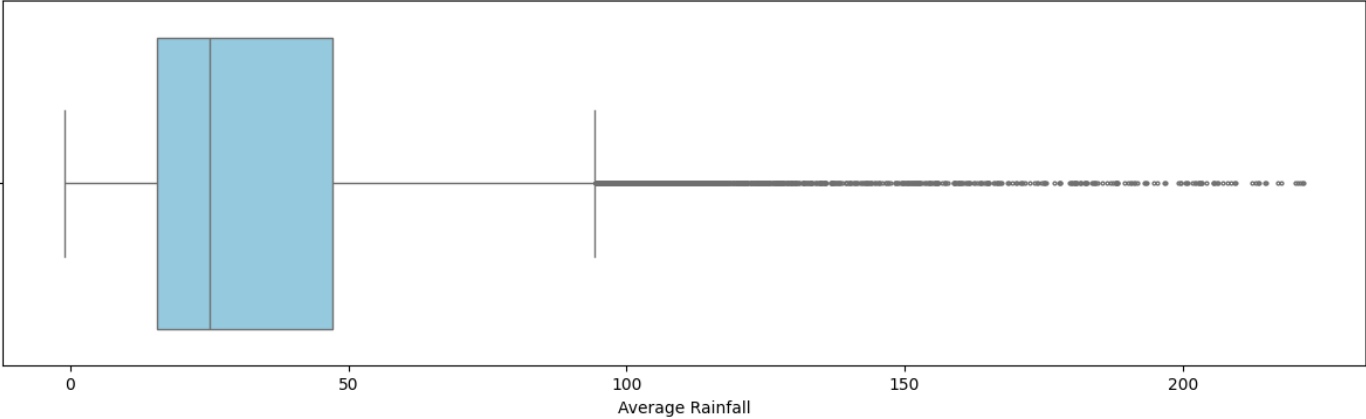
```

```

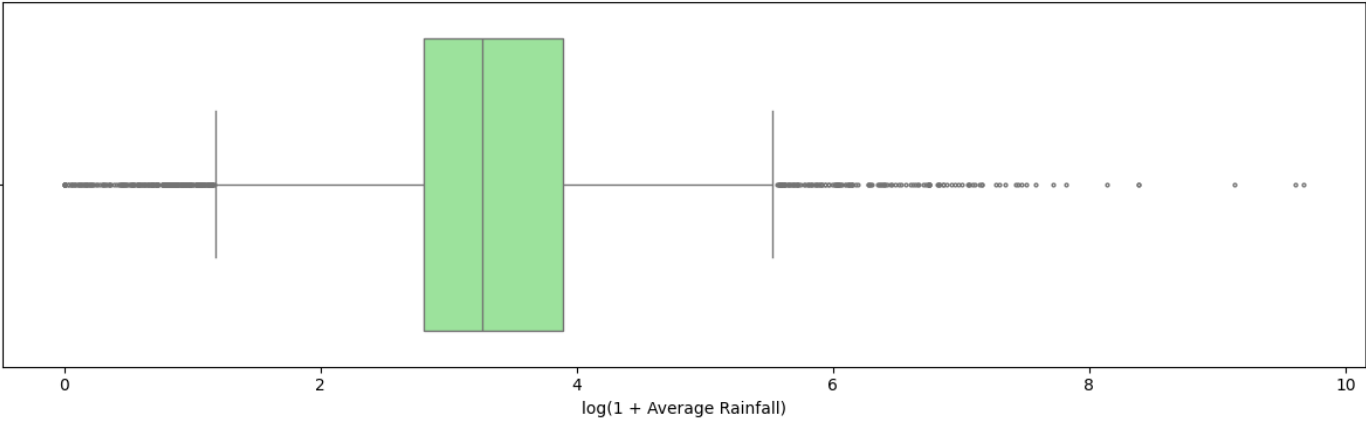
1 /usr/local/lib/python3.8/dist-packages/pandas/core/arraylike.py:402: RuntimeWarning: invalid value encountered in log1p
2     result = getattr(ufunc, method)(*inputs, **kwargs)

```

Distribution of Annual Average Rainfall (<99th Percentile)



Distribution of Annual Average Rainfall (Log Scale)



```
1 prcp_pdf[prcp_pdf.PRCP_yavg < 0]
```

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

year	COUNTRY_CODE	COUNTRY_NAME	PRCP_yavg
1874	UK	United Kingdom	-1.076712

Global Station Location View

```
1 prcp_pdf
```

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

year	COUNTRY_CODE	COUNTRY_NAME	PRCP_yavg
1874	UK	United Kingdom	-1.076712

2021	WQ	Wake Island [United States]	0.000000
2021	VE	Venezuela	0.000000
2021	AE	United Arab Emirates	0.000000
2021	BG	Bangladesh	0.000000
...	...	...	...
1951	None	None	4359.000000
2000	EK	Equatorial Guinea	4361.000000
1949	None	None	9271.000000
1950	None	None	14827.250000
1952	None	None	15875.000000

17731 rows × 4 columns

```
1 # =====
2 # Bubble map: PRCP stations per country
3 # (Spark aggregates → 218 rows → Pandas → plot)
4 # =====
5
6 import math
7 import pandas as pd
8 import geopandas as gpd
9 import matplotlib.pyplot as plt
10 from pyspark.sql import functions as F
11
12 # 1) Spark: 国家级聚合 (每站只计一次, 计算数量与国家“代表点”)
13 country_agg = (
14     prcp_stations
15     .select("ID", "COUNTRY_CODE", "LATITUDE", "LONGITUDE")
16     .dropna()
17     .filter((F.col("LATITUDE").between(-90, 90)) & (F.col("LONGITUDE").between(-180, 180)))
18     .dropDuplicates(["ID"]) # 避免同一站多次计入
19     .groupBy("COUNTRY_CODE")
20     .agg(
21         F.count("*").alias("station_count"),
22         F.avg("LATITUDE").alias("avg_lat"),
23         F.avg("LONGITUDE").alias("avg_lon"),
24     )
25 )
26
27 # 2) 收集到 Pandas (只有 ~218 行)
28 bubble_df = country_agg.toPandas()
```

```
1 # 3) 世界底图 (WGS84)
2 world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres")).to_crs(4326)
3
4 # 4) 把国家代表点做成 GeoDataFrame (WGS84)
5 bubbles = gpd.GeoDataFrame(
6     bubble_df,
7     geometry=gpd.points_from_xy(bubble_df["avg_lon"], bubble_df["avg_lat"]),
8     crs="EPSG:4326"
9 )
10
11 # 5) 可选: 投影到 Robinson (观感更好); 失败则保持 4326
12 target_crs = "ESRI:54030"
13 try:
14     world_proj = world.to_crs(target_crs)
15     bubbles_proj = bubbles.to_crs(target_crs)
16 except Exception:
17     target_crs = "EPSG:4326"
18     world_proj = world
19     bubbles_proj = bubbles
20
21 # 6) 气泡大小 (按平方根缩放, 保证视觉面积接近线性)
22 c = bubbles_proj["station_count"].astype(float)
23 c_min, c_max = float(c.min()), float(c.max())
24 # 你可以调 min/max_size 来控制视觉效果
25 min_size, max_size = 10, 800 # matplotlib scatter 的 s 是面积像素
26 if c_max > 0:
27     s = min_size + (max_size - min_size) * (c.pow(0.5) - math.sqrt(max(1.0, c_min))) / (math.sqrt(c_max) - math.sqrt(max(1.0, c_min)))
28     s = s.clip(lower=min_size)
29 else:
30     s = pd.Series([min_size] * len(c))
```

```

31
32 # 7) 绘图
33 fig, ax = plt.subplots(figsize=(16, 8))
34 world_proj.plot(ax=ax, color="#eef2f6", edgecolor="#c9d1da", linewidth=0.5, zorder=0)
35
36 # 国家代表点 (按站点数大小)
37 bubbles_proj.plot(
38     ax=ax,
39     markersize=s,      # 比例符号
40     alpha=0.65,
41     color="#1f77b4",
42     edgecolor="white",
43     linewidth=0.3,
44     zorder=2,
45 )
46
47 ax.set_axis_off()
48 ax.set_title(
49     f"Global PRCP Stations by Country (bubble size  $\propto$  station count)\nCRS: {target_crs}",
50     fontsize=13
51 )
52
53 # 8) 自定义气泡图例 (计算 5 个分级刻度 (对数刻度更均匀))
54 legend_vals = [50, 200, 1_000, 5_000, int(c_max)]
55
56 for legend_val in legend_vals:
57     if legend_val <= 0 or legend_val > c_max:
58         continue
59     # 用同样的缩放公式计算面积
60     area = min_size + (max_size - min_size) * (
61         (math.sqrt(legend_val) - math.sqrt(max(1.0, c_min)))
62         / (math.sqrt(c_max) - math.sqrt(max(1.0, c_min)))
63     )
64     plt.scatter([], [], s=area, color="#1f77b4", alpha=0.65,
65                 edgecolors="white", linewidths=0.3, label=f"{legend_val:,}")
66
67 leg = ax.legend(
68     scatterpoints=1,
69     frameon=True,
70     labelspacing=1.2,
71     title="Stations",
72     loc="lower left",
73     bbox_to_anchor=(0.02, 0.02),
74 )
75 leg.get_title().set_fontsize(10)
76
77 plt.tight_layout()
78 plt.savefig("./supplementary/global_prdp_stations_bubble.png", dpi=300, bbox_inches="tight")
79 plt.show()
80
81 # 9) 导出聚合表
82 out_csv = "./supplementary/prcp_station_country_bubbles.csv"
83 bubbles[["COUNTRY_CODE", "station_count", "avg_lat", "avg_lon"]].to_csv(out_csv, index=False)
84 print("Saved:", out_csv)

```

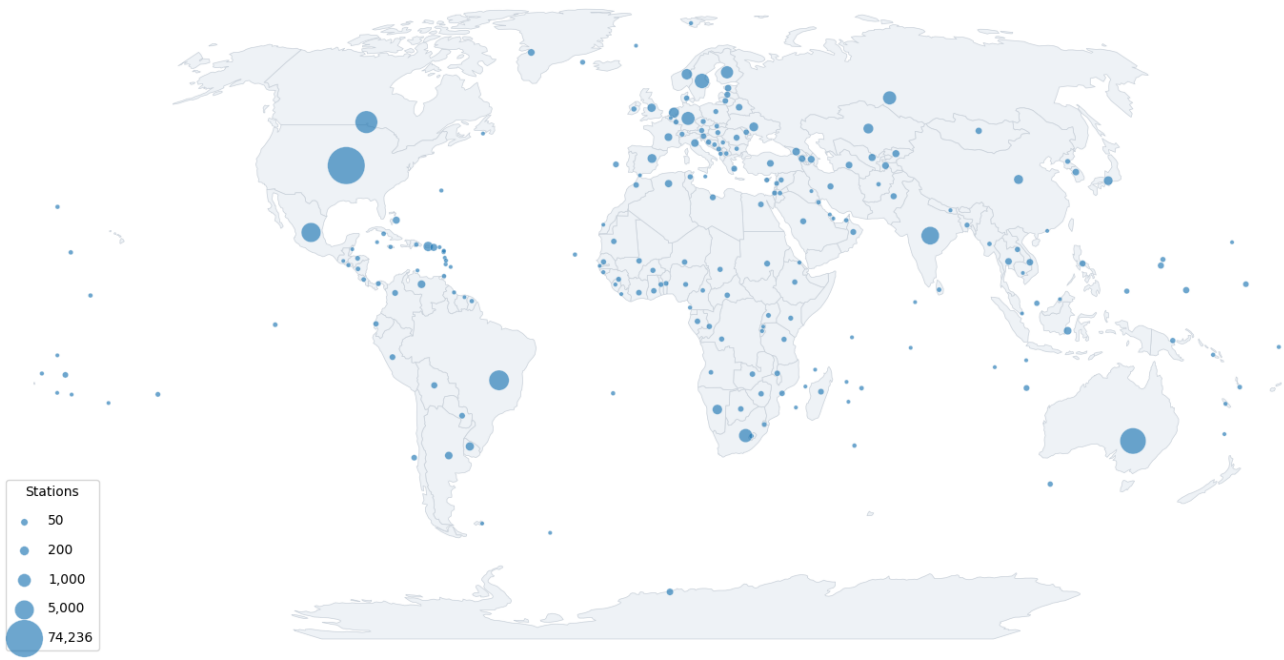
```

1 | /tmp/ipykernel_56/3098757670.py:2: FutureWarning: The geopandas.dataset module is deprecated and will be removed in GeoPandas 1.0. You can get
2 |   world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres")).to_crs(4326)

```



Global PRCP Stations by Country (bubble size  $\propto$  station count)  
CRS: ESRI:54030



1 | Saved: ./supplementary/prcp\_station\_country\_bubbles.csv

## PRCP Global Map

```
1 # =====
2 # 2024各国平均降水: FIPS10-6→ISO3映射 (自动补齐) + Choropleth
3 # 输出:
4 # 1) ./supplementary/fips_to_iso3_completed.csv (完整映射)
5 # 2) ./supplementary/precip_2024_by_country.csv (各国均值)
6 # 3) ./supplementary/precip_2024_choropleth.png (地图)
7 # =====
8
9 import os, math, json
10 import pandas as pd
11 import geopandas as gpd
12 import matplotlib.pyplot as plt
13 from pyspark.sql import functions as F
14
15 # ----- Spark 设置 (Arrow提速) -----
16 spark.conf.set("spark.sql.execution.arrow.pyspark.enabled", "true")
17
18 # ----- 0) 基础FIPS→ISO3字典 (常见+高频; 其余自动补全) -----
19 FIPS_TO_ISO3_BASE = {
20     # 常见主权国家 (示例, 非穷举; 其余会自动补齐)
21     "AF": "AFG", "AL": "ALB", "AG": "DZA", "AN": "AND", "AO": "AGO", "AC": "ATG", "AR": "ARG", "AM": "ARM",
22     "AS": "AUS", "AU": "AUT", "AJ": "AZE",
23     "BA": "BHS", "BG": "BGD", "BB": "BRB", "BO": "BEL", "BF": "BFA", "BU": "BGR", "BN": "BEN", "BD": "BMU",
24     "BT": "BTN", "BL": "BOL", "BR": "BIH", "BC": "BWA", "BR": "BRA", "BX": "BRN", "BUK": "MMR", # BU=保加利亚, BUK=缅甸 (旧称Burma)
25     "BY": "BLR", "BZ": "BLZ",
26     "CB": "KHM", "CM": "CMR", "CA": "CAN", "CV": "CPV", "CJ": "CYM", "CT": "CAF", "CD": "TCD", "CI": "CHL",
27     "CH": "CHN", "CO": "COL", "CN": "COM", "CF": "COG", "CW": "COG", "CS": "CRI", "IV": "CIV", "HR": "HRV",
28     "CU": "CUB", "CY": "CYP", "EZ": "CZE",
29     "DA": "DNK", "DJ": "DJI", "DO": "DMA", "DR": "DOM",
30     "EC": "ECU", "EG": "EGY", "ES": "SLV", "EK": "GNQ", "ER": "ERI", "EN": "EST", "ET": "ETH",
31     "FK": "FLK", "FJ": "FJI", "FI": "FIN", "FR": "FRA",
32     "GB": "GAB", "GA": "GMB", "GG": "GEO", "GM": "DEU", "GH": "GHA", "GI": "GIB", "GR": "GRC", "GL": "GRL",
33     "GJ": "GRD", "GT": "GTM", "GV": "GIN", "PU": "GNB", "GY": "GUY",
34     "HA": "HTI", "HO": "HND", "HU": "HUN",
35     "IC": "ISL", "IN": "IND", "ID": "IDN", "IR": "IRN", "IZ": "IRQ", "EI": "IRL", "IS": "ISR", "IT": "ITA",
36     "JM": "JAM", "JA": "JPN", "JO": "JOR",
37     "KZ": "KAZ", "KE": "KEN", "KR": "KIR", "KN": "KWT",
38     "LA": "LAO", "LG": "LVA", "LE": "LBN", "LI": "LSO", "LI": "LBR", "LY": "LBY", "LS": "LIE", "LH": "LTU",
39     "LU": "LUX",
40     "MK": "MDG", "MI": "MWI", "MY": "MYS", "MV": "MDV", "ML": "MLI", "MT": "MLT", "RM": "FSM", "MR": "MRT",
41     "MU": "MUS", "MX": "MEX", "MD": "MDA", "MN": "MCO", "MG": "MNG", "MJ": "MNE", "MO": "MAR", "MZ": "MOZ",
42     "WA": "NAM", "NR": "NRU", "NE": "NPL", "NL": "NLD", "NC": "NZL", "NU": "NIU", "NI": "NER", "NG": "NGA",
43     "NO": "NOR",
44     "MUQ": "OMN", "PK": "PAK", "PS": "PLW", "PM": "PAN", "PP": "PNG", "PA": "PRY", "PE": "PER", "RP": "PHL",
```

```

45 "PL": "POL", "PO": "PRT",
46 "QA": "QAT",
47 "RO": "ROU", "RS": "RUS", "RW": "RWA",
48 "SC": "KNA", "ST": "LCA", "VC": "VCT", "WS": "WSM", "SM": "SMR", "TP": "STP", "SA": "SAU", "SG": "SEN",
49 "RI": "SRB", "SE": "SYC", "SL": "SLE", "SN": "SGP", "LO": "SVK", "SI": "SVN", "BP": "SOM", "SF": "ZAF",
50 "OD": "SSD", "SP": "ESP", "CE": "LKA", "SU": "SDN", "NS": "SUR", "SV": "SWE", "SZ": "CHE", "SY": "SYR",
51 "TW": "TWN", "TI": "TJK", "TZ": "TZA", "TH": "THA", "TT": "TLS", "TO": "TGO", "TN": "TUN", "TU": "TUR",
52 "TX": "TKM", "TV": "TUV",
53 "UG": "UGA", "UP": "UKR", "AE": "ARE", "UK": "GBR", "US": "USA", "UY": "URY", "UZ": "UZB",
54 "NH": "VUT", "VE": "VEN", "VM": "VNM",
55 "YM": "YEM",
56 "ZA": "ZMB", "ZI": "ZWE",
57
58 # 常见属地/地区（部分库会并到主权国；此处给出常用ISO3便于merge）
59 "AQ": "ATA", # Antarctica
60 "BM": "MAC", # Macau（若你的数据用MO/MAC，按需调整）
61 "HK": "HKG",
62 "FO": "FRO", "GLP": "GLP", "GP": "GLP", "GF": "GUF", "PF": "PYF", "NCY": "NCL",
63 "RE": "REU", "PFY": "PYF", "GI": "GIB",
64 }

```

```

1 # ----- 1) 用Spark聚合 2024 年国家平均降水 -----
2 # 假设你已有 daily_prdp: 列 [ID, DATE(yyyy-mm-dd), PRCP_VALUE, year]
3 # 以及 prcp_stations: 列 [ID, COUNTRY_CODE, LATITUDE, LONGITUDE, ...]
4 # Define schema for Daily
5 stations_enriched = spark.read.parquet(stations_enriched_savepath)
6 daily_schema = StructType(
7     [
8         StructField("ID", StringType(), nullable=False),
9         StructField("DATE", StringType(), nullable=False),
10        StructField("ELEMENT", StringType(), nullable=False),
11        StructField("VALUE", FloatType(), nullable=False),
12        StructField("MEASUREMENT_FLAG", StringType(), nullable=True),
13        StructField("QUALITY_FLAG", StringType(), nullable=True),
14        StructField("SOURCE_FLAG", StringType(), nullable=True),
15        StructField("OBSERVATION_TIME", StringType(), nullable=True),
16    ]
17 )
18
19 # load daily and check daily schema for later join parameter on = ""
20 daily = spark.read.csv(paths["daily"], schema=daily_schema)
21
22 daily_prdp = (
23     daily.filter(F.col("ELEMENT") == "PRCP")
24     .select(["ID", "DATE", "VALUE"])
25     .withColumn("DATE", F.to_date(F.col("DATE"), "yyyyMMdd"))
26     .withColumnRenamed("VALUE", "PRCP_VALUE")
27     .withColumn("year", F.year("DATE"))
28 )
29
30 prcp_stations = daily_prdp.join(
31     F.broadcast(stations_enriched.filter(F.col("PRCP") == "1")),
32     on="ID",
33     how="left",
34 ).withColumnRenamed("NAME", "STATION_NAME")
35
36 prcp_2024 = (
37     daily_prdp
38     .filter(F.col("year") == 2024)
39     .select("ID", "PRCP_VALUE")
40     .groupBy("ID")
41     .agg(F.avg("PRCP_VALUE").alias("PRCP_2024_station_avg"))
42 )
43
44 # 连接国家代码（只需要ID-COUNTRY_CODE）
45 station_country = prcp_stations.select("ID", "COUNTRY_CODE").dropna().dropDuplicates(["ID"])
46
47 country_avg_2024 = (
48     prcp_2024.join(station_country, on="ID", how="inner")
49     .groupBy("COUNTRY_CODE")
50     .agg(F.avg("PRCP_2024_station_avg").alias("PRCP_2024_country_avg"),
51         F.countDistinct("ID").alias("n_stations_2024"),
52         F.avg("PRCP_2024_station_avg").alias("check_same")) # 冗余校验
53 )
54
55 # 拉到 pandas (只有 ~218 行)
56 avg_pdf = country_avg_2024.toPandas()

```

```

1 # ----- 2) 先用基础字典做FIPS-ISO3映射 -----

```

```

2 avg_pdf["iso_a3"] = avg_pdf["COUNTRY_CODE"].map(FIPS_TO_ISO3_BASE)
3
4 # ----- 3) 对未映射成功的FIPS, 做“轻量空间补齐” (每FIPS取站点平均坐标~落国界) -----
5 need_fill = avg_pdf[avg_pdf["iso_a3"].isna()][["COUNTRY_CODE"]].unique().tolist()
6
7 if len(need_fill) > 0:
8     # 为这些FIPS计算代表点 (该FIPS内所有站点的经纬度均值)
9     fips_reps = (
10         prcp_stations
11         .select("COUNTRY_CODE", "LATITUDE", "LONGITUDE")
12         .where(F.col("COUNTRY_CODE").isin(need_fill))
13         .groupBy("COUNTRY_CODE")
14         .agg(F.avg("LATITUDE").alias("avg_lat"), F.avg("LONGITUDE").alias("avg_lon"))
15         .toPandas()
16     )
17
18     # 点→GeoDataFrame (WGS84)
19     reps_gdf = gpd.GeoDataFrame(
20         fips_reps,
21         geometry=gpd.points_from_xy(fips_reps["avg_lon"], fips_reps["avg_lat"]),
22         crs="EPSG:4326"
23     )
24     # 世界底图 (WGS84)
25     world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres")).to_crs(4326)
26     world = world[~world["iso_a3"].isin(["-99"])[["iso_a3", "name", "geometry"]]
27
28     # 空间匹配 (点落国家多边形)
29     reps_joined = gpd.sjoin(reps_gdf, world, how="left", predicate="within")
30     # 生成补齐映射
31     filled_map = dict(zip(reps_joined["COUNTRY_CODE"], reps_joined["iso_a3"]))
32
33     # 合并到基础字典
34     FIPS_TO_ISO3_COMPLETED = {**FIPS_TO_ISO3_BASE, **filled_map}
35 else:
36     FIPS_TO_ISO3_COMPLETED = FIPS_TO_ISO3_BASE.copy()
37
38 # 应用完整映射
39 avg_pdf["iso_a3"] = avg_pdf["iso_a3"].fillna(avg_pdf["COUNTRY_CODE"].map(FIPS_TO_ISO3_COMPLETED))
40
41 # 落盘保存完整映射 (便于复现)
42 os.makedirs("./supplementary", exist_ok=True)
43 pd.Series(FIPS_TO_ISO3_COMPLETED).rename("ISO3").to_csv(
44     "./supplementary/fips_to_iso3_completed.csv", header=True
45 )
46
47 # ----- 4) 与世界底图合并并绘图 (No data 灰色) -----
48 world = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres")).to_crs(4326)
49 world = world[~world["iso_a3"].isin(["-99"])]
50
51 # 合并 (左表为世界底图, 确保无观测国家也在图上)
52 world_join = world.merge(
53     avg_pdf[["iso_a3", "PRCP_2024_country_avg", "n_stations_2024"]],
54     on="iso_a3", how="left"
55 )
56
57 # 导出表
58 world_join[["iso_a3", "name", "PRCP_2024_country_avg", "n_stations_2024"]]\
59     .to_csv("./supplementary/precip_2024_by_country.csv", index=False)
60
61 # 选择投影 (Robinson观感好; 失败则保持WGS84)
62 target_crs = "ESRI:54030"
63 try:
64     world_plot = world_join.to_crs(target_crs)
65 except Exception:
66     target_crs = "EPSG:4326"
67     world_plot = world_join

```

```

1 # 绘图 (无观测国家灰色)
2 fig, ax = plt.subplots(figsize=(16, 8))
3 world_plot.plot(
4     ax=ax,
5     column="PRCP_2024_country_avg",
6     cmap="YlGnBu",
7     legend=True,
8     legend_kwds={"label": "Average Rainfall in 2024 (country mean of station means)"},
9     edgecolor="#d3d9df",
10    linewidth=0.4,
11    missing_kwds={"color": "#f0f0f0", "edgecolor": "#d3d9df", "hatch": "///", "label": "No data"},
12 )
13 # ax.set_axis_off()
14 ax.set_title(
15     "Average Rainfall by Country (2024)\n"

```

```

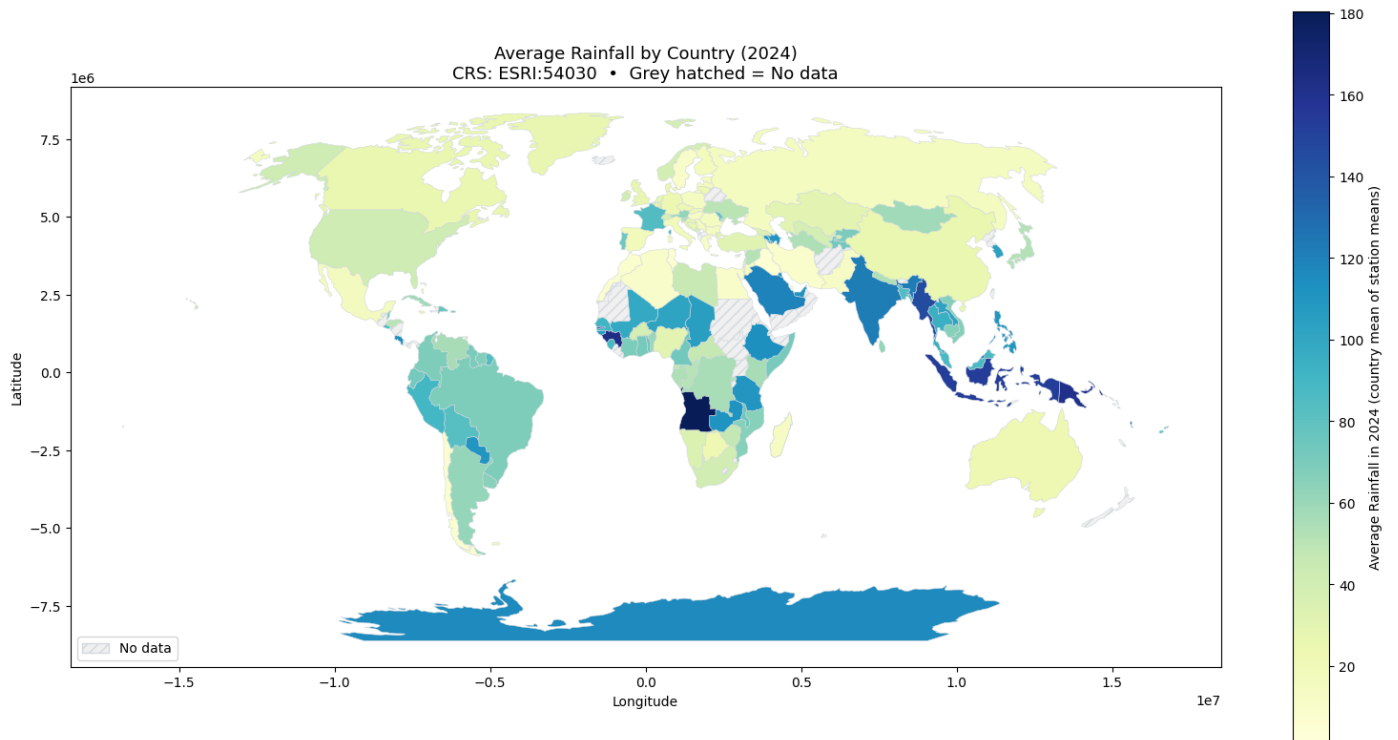
16     f"CRS: {target_crs} • Grey hatched = No data",
17     fontsize=13
18 )
19
20 # display longitude and latitude axis
21 ax.set_xlabel("Longitude")
22 ax.set_ylabel("Latitude")
23
24 # 给“无数据”加图例项 (有的版本会自动, 有的需要手动补)
25 handles, labels = ax.get_legend_handles_labels()
26 if "No data" not in labels:
27     from matplotlib.patches import Patch
28     handles.append(Patch(facecolor="#f0f0f0", edgecolor="#d3d3d3", hatch="///", label="No data"))
29     ax.legend(handles=handles, loc="lower left")
30
31 plt.tight_layout()
32 plt.savefig("./supplementary/precip_2024_choropleth.png", dpi=220, bbox_inches="tight")
33 plt.show()
34
35 print("✓ Saved mapping to ./supplementary/fips_to_iso3_completed.csv")
36 print("✓ Saved country table to ./supplementary/precip_2024_by_country.csv")
37 print("✓ Saved map to ./supplementary/precip_2024_choropleth.png")

```

```

1 /tmp/ipykernel_44/1411382327.py:25: UserWarning: Legend does not support handles for PatchCollection instances.
2 See: https://matplotlib.org/stable/tutorials/intermediate/legend_guide.html#implementing-a-custom-legend-handler
3     handles, labels = ax.get_legend_handles_labels()

```



```

1 ✓ Saved mapping to ./supplementary/fips_to_iso3_completed.csv
2 ✓ Saved country table to ./supplementary/precip_2024_by_country.csv
3 ✓ Saved map to ./supplementary/precip_2024_choropleth.png

```

```

1 import pandas as pd
2 import geopandas as gpd
3 import matplotlib.pyplot as plt
4 from matplotlib.ticker import MultipleLocator, FuncFormatter
5
6 # 1) 读入几何和csv (无需再计算)
7 world_base = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres"))[
8     ["iso_a3", "name", "geometry"]
9 ]
10 prec = pd.read_csv("./supplementary/precip_2024_by_country.csv")
11
12 # 如果你的csv列名不同, 请在这里对齐列名
13 # 例如: prec.rename(columns={"iso3": "iso_a3"}, inplace=True)

```

```

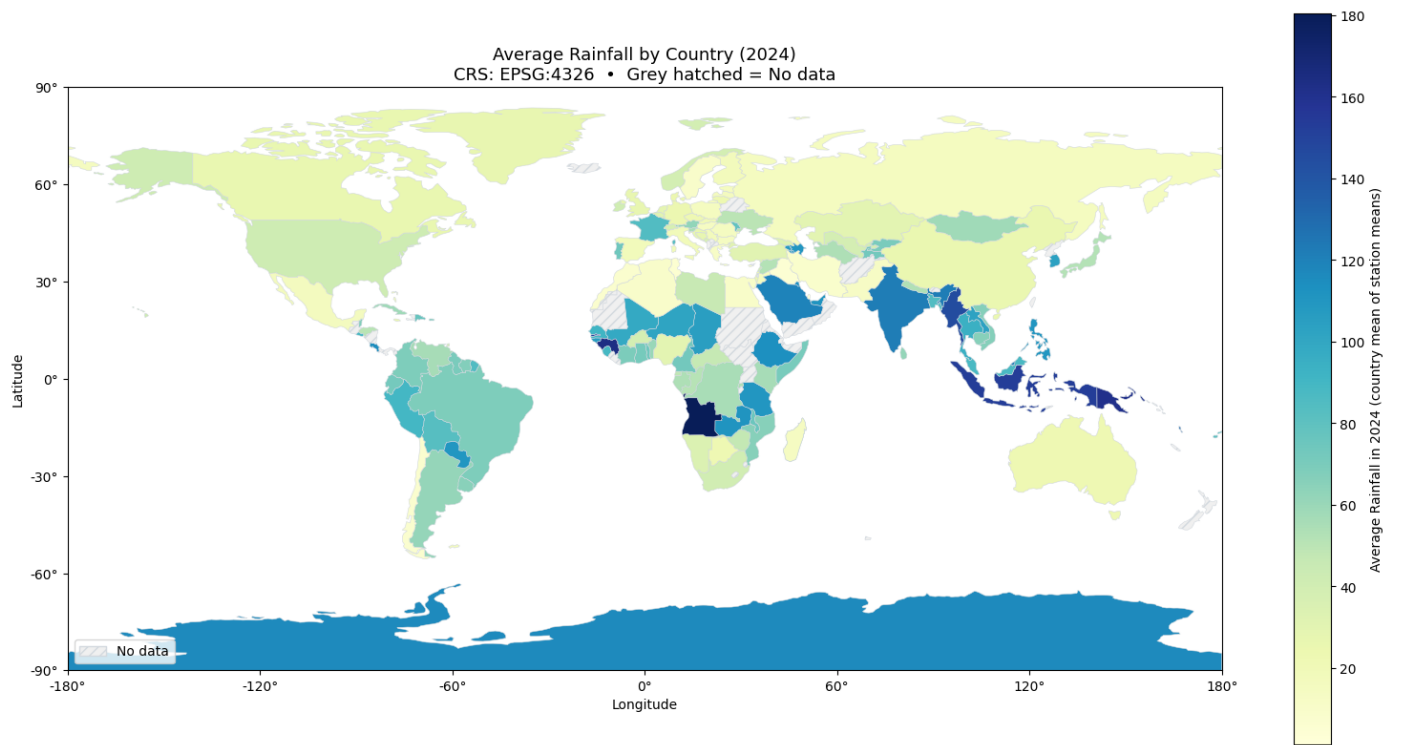
14
15 # 2) 合并到GeoDataFrame
16 world_plot = world_base.merge(
17     prec[["iso_a3", "PRCP_2024_country_avg"]],
18     on="iso_a3", how="left"
19 ).to_crs("EPSG:4326")
20
21 # 3) 画图: 坐标轴显示度数
22 fig, ax = plt.subplots(figsize=(16, 8))
23 world_plot.plot(
24     ax=ax,
25     column="PRCP_2024_country_avg",
26     cmap="YlGnBu",
27     legend=True,
28     legend_kwds={"label": "Average Rainfall in 2024 (country mean of station means)"},
29     edgecolor="#d3d9df", linewidth=0.4,
30     missing_kwds={"color": "#f0f0f0", "edgecolor": "#d3d9df", "hatch": "///", "label": "No data"}
31 )
32
33 ax.set_xlabel("Longitude")
34 ax.set_ylabel("Latitude")
35
36 def deg(x, pos): # 度数格式
37     return f"{int(x)}°"
38
39 ax.xaxis.set_major_locator(MultipleLocator(60))
40 ax.yaxis.set_major_locator(MultipleLocator(30))
41 ax.xaxis.set_major_formatter(FuncFormatter(deg))
42 ax.yaxis.set_major_formatter(FuncFormatter(deg))
43 ax.set_xlim(-180, 180)
44 ax.set_ylim(-90, 90)
45
46 ax.set_title(
47     "Average Rainfall by Country (2024)\n"
48     "CRS: EPSG:4326 • Grey hatched = No data",
49     fontsize=13
50 )
51
52 # 补上"No data"图例 (如果自动没加上)
53 handles, labels = ax.get_legend_handles_labels()
54 if "No data" not in labels:
55     from matplotlib.patches import Patch
56     handles.append(Patch(facecolor="#f0f0f0", edgecolor="#d3d9df", hatch="///", label="No data"))
57     ax.legend(handles=handles, loc="lower left")
58
59 plt.tight_layout()
60 plt.savefig("./supplementary/precip_2024_choropleth_from_csv_wgs84.png", dpi=220, bbox_inches="tight")
61 plt.show()

```

```

1 /tmp/ipykernel_44/2271823713.py:7: FutureWarning: The geopandas.dataset module is deprecated and will be removed in GeoPandas 1.0. You can get
2   world_base = gpd.read_file(gpd.datasets.get_path("naturalearth_lowres"))[
3 /tmp/ipykernel_44/2271823713.py:53: UserWarning: Legend does not support handles for PatchCollection instances.
4 See: https://matplotlib.org/stable/tutorials/intermediate/legend_guide.html#implementing-a-custom-legend-handler
5   handles, labels = ax.get_legend_handles_labels()

```



```
1 | stop_spark()
```

```
1 | 25/09/14 13:03:42 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](#)

```
1 |
```