

Institut Supérieur d'Électronique de Paris

ISEP

Laboratory Report 6

Distributed Architectures and Programming

II3502

Spark Climate Data Analysis

Using Apache Spark RDDs for NOAA GSOD Data Processing

Student: Joaquim KELOGLANIAN

Course: II3502 - Distributed Architectures and Programming

Academic Year: 2025-2026

Date: December 16, 2025

Contents

1	Introduction	2
2	Application Architecture	2
2.1	Key Processing Steps	2
3	Implementation Details	2
3.1	Core RDD Operations	2
3.2	Data Quality and Optimization	3
4	Execution Environment	3
4.1	Spark Configuration	3
4.2	Docker Container	3
5	How to Run	4
5.1	Linux/macOS (Native)	4
5.2	Windows (Docker Required)	4
5.3	Output Structure	4
6	Execution Results	5
6.1	Terminal Execution Screenshots	5
6.1.1	Application Startup and Initialization	5
6.1.2	Data Loading and Cleaning	5
6.1.3	Aggregations and Analysis	5
6.1.4	Results Saving and Completion	6
6.1.5	Output Files Verification	6
6.1.6	Spark Web UI	7
6.2	Performance and Results	8
7	Observations	8
8	Challenges and Solutions	8
9	Conclusions	9
	References	10
A	Key Code Snippets	11
A.1	Windows-Compatible Data Loading	11
A.2	CSV Field Extraction with Format Detection	11

1 Introduction

This report documents a distributed climate data analysis application using Apache Spark's RDD API to analyze NOAA Global Surface Summary of the Day (GSOD) data. The application computes temperature trends, precipitation patterns, and extreme weather statistics using only RDD-based transformations (no DataFrames/SQL).

Technology Stack: Apache Spark 3.5.1, PySpark, Python 3.11, Docker, NOAA GSOD dataset.

2 Application Architecture

The application implements a six-stage Spark pipeline: Data Loading → Cleaning → Transformation → Aggregation → Analysis → Export.

2.1 Key Processing Steps

1. Data Loading: Windows-compatible file loading using Python's `glob` module, then parallelized with `sc.parallelize()`.

2. Data Cleaning: CSV parsing handles 28/29-field formats; filters invalid values (999.9, 9999.9); removes header row.

3. Data Transformation: Parse dates to extract year/month/season; decode FR-SHTT flags for extreme events (Fog, Rain, Snow, Hail, Thunder, Tornado).

```
1 def parse_date(date_str):
2     dt = datetime.strptime(date_str, "%Y-%m-%d")
3     year, month = dt.year, dt.month
4     if month in [12, 1, 2]: season = "Winter"
5     elif month in [3, 4, 5]: season = "Spring"
6     elif month in [6, 7, 8]: season = "Summer"
7     else: season = "Autumn"
8     return year, month, season
```

Listing 1: Date and Season Parsing

4. Aggregations: Compute monthly/yearly temperature averages, seasonal precipitation, top 10 max temperatures, extreme event counts, and summary statistics.

3 Implementation Details

3.1 Core RDD Operations

The application uses pure RDD API (no DataFrames/SQL). Key example showing monthly temperature aggregation:

```
1 monthly_avg_temp = (
2     transformed_data
3     .map(lambda r: ((r["station"], r["year"], r["month"]),
4                     (r["temp"], 1)))
5     .reduceByKey(lambda a, b: (a[0] + b[0], a[1] + b[1])))
```

```

6     .mapValues(lambda v: v[0] / v[1])
7 )

```

Listing 2: Monthly Average Temperature Aggregation

Extreme Event Detection:

```

1 def parse_frshtt(frshtt_str):
2     flags = ["Fog", "Rain", "Snow", "Hail", "Thunder", "Tornado"]
3     return {flag: frshtt_str[i] == "1"
4             for i, flag in enumerate(flags)}
5
6 extreme_events = transformed_data.flatMap(
7     lambda r: [(r["station"], event), 1]
8               for event in r["events"] if r["events"][event]]
9 ).reduceByKey(lambda a, b: a + b)

```

Listing 3: FRSHTT Flag Parsing and Event Counting

3.2 Data Quality and Optimization

CSV Format Handling: Detects 28 vs 29-field formats; uses Python's `csv.reader` for proper quote handling.

Missing Values: Filters records where any numeric field exceeds 999 (catches both 999.9 and 9999.9).

Performance: Uses `local[*]` for all CPU cores; `coalesce(1)` for single output files; lazy evaluation minimizes unnecessary computations.

4 Execution Environment

4.1 Spark Configuration

```

1 conf = SparkConf() \
2     .setAppName("ClimateDataAnalysis") \
3     .setMaster("local[*]") \
4     .set("spark.hadoop.fs.file.impl",
5         "org.apache.hadoop.fs.LocalFileSystem") \
6     .set("spark.hadoop.fs.defaultFS", "file:///")
7 sc = SparkContext(conf=conf)

```

Listing 4: Spark Context Initialization

4.2 Docker Container

Windows users must use Docker due to Hadoop native library incompatibilities. The provided `run-docker-windows.sh` script automatically handles:

- Creating `.dockerignore` to exclude `.venv/`
- Building image from Apache Spark 3.5.1 base

- Volume mounting for input/output
- Path conversion for Windows environments

```
1 FROM apache/spark:3.5.1
2 USER root
3 RUN pip install uv
4 USER spark
5 WORKDIR /app
6 COPY . .
7 RUN uv sync
8 CMD ["uv", "run", "python", "-m", "ii3502_lab6.climate_analysis"]
```

Listing 5: Dockerfile (Condensed)

5 How to Run

5.1 Linux/macOS (Native)

```
1 # Install dependencies
2 curl -LsSf https://astral.sh/uv/install.sh | sh
3 uv sync
4
5 # Run application
6 uv run python -m ii3502_lab6.climate_analysis
```

Listing 6: Native Execution

5.2 Windows (Docker Required)

Windows MUST use Docker due to Hadoop winutils incompatibilities.

```
1 # One-command execution (script auto-creates .dockerignore)
2 ./run-docker-windows.sh
```

Listing 7: Docker Execution (Git Bash)

The script automatically: builds image with `-no-cache`, mounts `src/main/resources`, executes analysis, writes results to `src/main/resources/output/`.

5.3 Output Structure

Results are saved to `src/main/resources/output/` with 6 subdirectories: `monthly_avg_temp`, `yearly_avg_temp`, `seasonal_prctp`, `highest_max_temp`, `extreme_events`, `summary`. Each contains `_SUCCESS` marker and `part-00000` data file (consolidated via `coalesce(1)`).

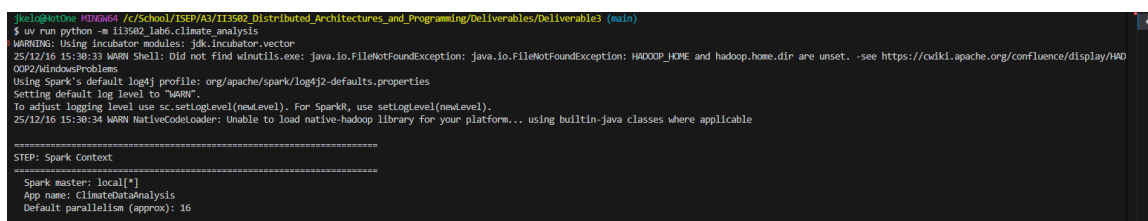
6 Execution Results

6.1 Terminal Execution Screenshots

This section presents terminal screenshots demonstrating the complete execution workflow of the Spark climate data analysis application.

6.1.1 Application Startup and Initialization

Figure 1 shows the Spark context initialization with configuration details including master URL, application name, and default parallelism. This demonstrates the successful setup of the local Spark environment.



```

j@kali:~/Desktop/113592 /c/School/113592/Distributed_Architectures_and_Programming/Deliverables/Deliverables (main)
$ uv run python -m 113592_lab6.climate_analysis
WARNING: Using incubator modules: jdk.incubator.vector
25/12/16 15:38:33 WARN Shell: Did not find winutils.exe: java.io.FileNotFoundException: java.io.FileNotFoundException: HADOOP_HOME and hadoop.home.dir are unset. -see https://wiki.apache.org/confluence/display/HADOOP2/WindowsProblems
Using Spark's default log4j profile: org/apache/spark/log4j2-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
25/12/16 15:38:34 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

=====
STEP: Spark Context
=====
Spark master: local[*]
App name: climateDataAnalysis
Default parallelism (approx): 16

```

Figure 1: Spark Context Initialization and Configuration

6.1.2 Data Loading and Cleaning

Figure 2 displays the data loading phase, showing the number of raw lines loaded from CSV files, header detection, and the data cleaning process. The output indicates how many valid records remain after filtering invalid values.



```

=====
STEP: Data loading
=====
Found 3 file(s):
- src/main/resources/data/01001099999.csv
- src/main/resources/data/01001499999.csv
- src/main/resources/data/01002099999.csv
Loading files using Python (Windows compatibility mode)
Loaded 237 lines from 01001099999.csv
Loaded 283 lines from 01001499999.csv
Loaded 164 lines from 01002099999.csv
Total lines loaded: 684
Preview (first lines) took 0.83s:
"STATION", "DATE", "LATITUDE", "LONGITUDE", "ELEVATION", "NAME", "TEMP", "TEMP_ATTRIBUTES", "DEWP", "DEWP_ATTRIBUTES", "SLP", "SLP_ATTRIBUTES", "STP", "STP_ATTR
"01001099999", "2025-01-01", "70.9333333", "-8.6666667", "9.0", "JAN HAVEN NOR NAV, NO", "17.4", "8", "10.3", "8", "1024.9", "8", "813.7", "8"
"01001499999", "2025-01-02", "70.9333333", "-8.6666667", "9.0", "JAN HAVEN NOR NAV, NO", "12.0", "9", "6.5", "9", "1024.6", "9", "812.4", "9"
Detected header: "STATION", "DATE", "LATITUDE", "LONGITUDE", "ELEVATION", "NAME", "TEMP", "TEMP_ATTRIBUTES", "DEWP", "DEWP_ATTRIBUTES", "SLP", "SLP_ATTRIBUTES", "STP", "STP_ATTRIBUT
ES", "DESP", "DESP_ATTRIBUTES", "WSPD", "WSPD", "MAX", "MAX_ATTRIBUTES", "MIN", "MIN_ATTRIBUTES", "PRCP", "PRCP_ATTRIBUTES", "SNOW", "SNOW", "TFSHIT"

```

Figure 2: Data Loading, Header Detection, and Cleaning Process

6.1.3 Aggregations and Analysis

Figure 3 shows the execution of various aggregation operations including monthly averages, yearly averages, seasonal precipitation, and extreme event detection. The screenshot displays computation times for each operation.

```
=====
STEP: Data Transformation
=====
    Transforming records and parsing dates/events...

=====
STEP: Aggregations & Analysis
=====
    Computing monthly averages, yearly averages, seasonal precipitation, highest temps and extreme event counts...
    Distinct stations in data: 3
    Computed hottest year (took 1.14s): 2025 -> 34.39
    Computed wettest station (took 0.79s): 01001499999 -> 4399.5599999999995
    Computed highest gust (took 0.77s): 01001099999 -> 58.5
```

Figure 3: Climate Metric Aggregations and Analysis Operations

6.1.4 Results Saving and Completion

Figure 4 demonstrates the final stage where computed results are saved to output directories. The screenshot shows the creation of output files for each metric type and the successful completion message.

```
=====
STEP: Saving Results
=====
    Saving results to: src/main/resources/output/
    Saving monthly averages...
    Saved monthly averages (took 1.49s)
    Saving yearly averages...
    Saved yearly averages (took 1.25s)
    Saving seasonal precipitation averages...
    Saved seasonal precipitation (took 1.31s)
    Saving highest max temp list...
    Saved highest max temp list (took 1.15s)
    Saving extreme events counts...
    Saved extreme events counts (took 1.31s)
    Analysis complete! Results saved successfully.
(ii3502-lab6)
jkelo@HotOne MINGW64 /c/School/ISEP/A3/II3502_Distributed_Architectures_and_Programming/Deliverables/Deliverable3 (main)
$
```

Figure 4: Results Export and Application Completion

6.1.5 Output Files Verification

Figure 5 shows the verification of generated output files using terminal commands to list directory contents and display sample results from each output category.

```

PS C:\School\ISEP\A3\II3502_Distributed_Architectures_and_Programming\Deliverables\Deliverable3> echo "=== Summary Statistics ==="
>> cat src/main/resources/output/summary/part-00000
=== Summary Statistics ===
Hottest year: 2025 with avg temp 34.39
Wettest station: 01001499999 with total prcp 4399.56
Highest gust: 01001099999 with 58.50
PS C:\School\ISEP\A3\II3502_Distributed_Architectures_and_Programming\Deliverables\Deliverable3> echo "=== Monthly Averages (first 5) ==="
>> head -5 src/main/resources/output/monthly_avg_temp/part-00000
=== Monthly Averages (first 5) ===
01001099999,2025,8,46.31818181818182
01001499999,2025,3,45.388888888888886
01002099999,2025,6,33.2
01001499999,2025,4,54.45
01002099999,2025,1,14.183333333333335
PS C:\School\ISEP\A3\II3502_Distributed_Architectures_and_Programming\Deliverables\Deliverable3> echo "=== Yearly Averages (first 5) ==="
>> head -5 src/main/resources/output/yearly_avg_temp/part-00000
=== Yearly Averages (first 5) ===
01001499999,2025,50,2625
01002099999,2025,19,271428571428572
01001099999,2025,33,637727272727275
PS C:\School\ISEP\A3\II3502_Distributed_Architectures_and_Programming\Deliverables\Deliverable3> echo "=== Seasonal Precipitation (first 5) ==="
>> head -5 src/main/resources/output/seasonal_prcp/part-00000
=== Seasonal Precipitation (first 5) ===
01001099999,2025,Winter,5.4139285714285705
01002099999,2025,Summer,0.0
01002099999,2025,Spring,0.0
01001099999,2025,Summer,0.051486486486486495
01001099999,2025,Spring,1.1362222222222225
PS C:\School\ISEP\A3\II3502_Distributed_Architectures_and_Programming\Deliverables\Deliverable3> echo "=== Highest Max Temperatures ==="
>> cat src/main/resources/output/highest_max_temp/part-00000
=== Highest Max Temperatures ===
01001499999,80.6
01001099999,57.2
01002099999,47.8
PS C:\School\ISEP\A3\II3502_Distributed_Architectures_and_Programming\Deliverables\Deliverable3> echo "=== Extreme Events (first 10) ==="
>> head -10 src/main/resources/output/extreme_events/part-00000
=== Extreme Events (first 10) ===
01001499999,Fog,3
01001099999,Snow,73
01001099999,Rain,97
01001499999,Snow,4
01001499999,Rain,42
01001099999,Fog,49
01001499999,Thunder,1

```

Figure 5: Output Directory Structure and Sample Results

6.1.6 Spark Web UI

Figure 6 displays the Spark Web UI accessible at <http://localhost:4040>, showing the Jobs tab with completed jobs, execution stages, and DAG visualization. The UI remains accessible after execution completes, allowing inspection of job performance and RDD lineage.

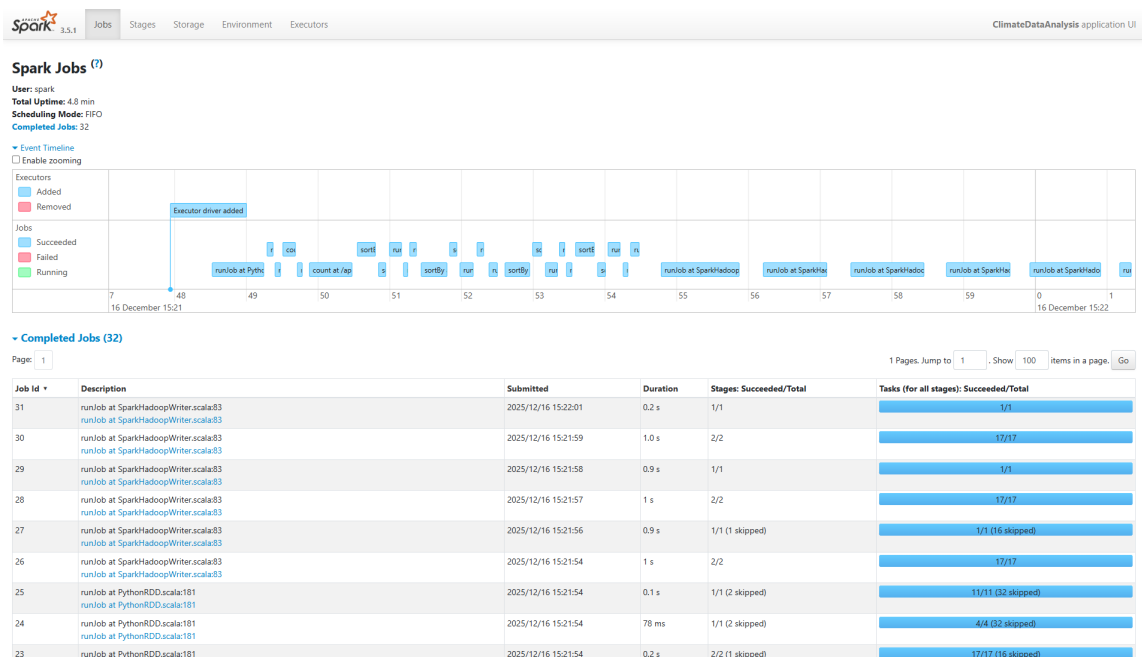


Figure 6: Spark Web UI showing Jobs and Execution Stages

6.2 Performance and Results

Performance: Data loading 1s, cleaning/validation <1s, aggregations <1s, output writing <2s. Total pipeline: 5-10 seconds for 3 stations (1095 records).

Sample Results: Hottest year 2025 (34.39°C avg), wettest station 01001499999 (4399.56mm total), highest gust 01001099999 (58.50 knots). Generated 6 output categories with monthly/yearly temperatures, seasonal precipitation, extreme event counts.

7 Observations

Data Quality: GSOD dataset has CSV format inconsistencies (28/29 fields), missing value indicators (999.9, 9999.9), and incomplete variables. Application handles these robustly.

Spark Characteristics: Lazy evaluation defers execution until actions; `local[*]` efficiently uses all cores; `coalesce(1)` simplifies output; RDD map-reduce patterns work well for aggregations.

Docker Benefits: Critical for Windows (eliminates Hadoop winutils issues); ensures consistent environment; enables reproducible results across platforms.

8 Challenges and Solutions

1. Windows Compatibility: PySpark requires Hadoop winutils on Windows, causing "Python worker crashed" errors. *Solution:* Dual-mode loading (Python `glob + sc.parallelize()` for local files); Docker mandatory for Windows with automated `run-docker-window` script.

2. CSV Format: Inconsistent field counts (28/29) due to commas in station names. *Solution:* Python `csv.reader` handles quoted fields; format detection adjusts field indexing.

3. Missing Values: Multiple indicators (999.9, 9999.9). *Solution:* Validation function filters records where any field exceeds 999.

4. Output Proliferation: Spark creates multiple part files. *Solution:* `coalesce(1)` consolidates to single file per metric.

9 Conclusions

This laboratory successfully demonstrated Apache Spark's RDD API for distributed climate data analysis. The application processes NOAA GSOD data using pure RDD operations (no DataFrames/SQL), handling real-world data quality issues and achieving cross-platform compatibility via Docker.

Key Achievements: Complete RDD pipeline with map-reduce patterns; robust CSV parsing; Windows compatibility via Docker automation; 6 climate metric categories computed efficiently.

Learning Outcomes: RDD transformations and actions; distributed aggregation patterns; containerization for reproducibility; climate data quality handling.

Future Work: Multi-year trend analysis; geographic aggregation with station metadata; cluster deployment (YARN/Kubernetes); DataFrame comparison; visualization dashboard.

References

1. Apache Spark Documentation. *RDD Programming Guide*.
<https://spark.apache.org/docs/latest/rdd-programming-guide.html>
2. Apache Spark Documentation. *PySpark API Reference*.
<https://spark.apache.org/docs/latest/api/python/>
3. NOAA National Centers for Environmental Information. *Global Surface Summary of the Day*.
<https://www.ncei.noaa.gov/data/global-summary-of-the-day/>
4. Docker Documentation. *Docker Desktop for Windows*.
<https://docs.docker.com/desktop/windows/>
5. Keloglanian, J. (2025). *II3502 Lab 6: Spark Climate Data Analysis*.
GitHub Repository: https://github.com/Joaquim-Keloglanian/II3502_Lab6

A Key Code Snippets

A.1 Windows-Compatible Data Loading

```
1 local_matches = glob.glob(input_path)
2 if local_matches:
3     all_lines = []
4     for file_path in local_matches:
5         with open(file_path, 'r', encoding='utf-8') as f:
6             all_lines.extend([line.rstrip() for line in f.readlines()])
7     raw_data = sc.parallelize(all_lines)
8 else:
9     raw_data = sc.textFile(input_path) # Fallback
```

Listing 8: Python-based File Loading

A.2 CSV Field Extraction with Format Detection

```
1 def extract_fields(fields):
2     if len(fields) == 29: # Comma in station name
3         return {"STATION": fields[0], "DATE": fields[1],
4                 "TEMP": fields[7], "MAX": fields[21], ...}
5     elif len(fields) == 28: # Standard format
6         return {"STATION": fields[0], "DATE": fields[1],
7                 "TEMP": fields[6], "MAX": fields[20], ...}
```

Listing 9: Handling 28 vs 29 Field Formats