



# CIS 4560 Term Project Tutorial



**Authors:** Kathlyne Alilain, Jennifer Mazas

**Instructor:** [Jongwook Woo](#)

**Date:** 05/05/2025

## Lab Tutorial

Kathlyne Alilain ([kalilai2@calstatela.edu](mailto:kalilai2@calstatela.edu))

Jennifer Mazas([jmazas@calstatela.edu](mailto:jmazas@calstatela.edu))

05/05/2025

## Evolution of Europe's Climate Monitoring Infrastructure

---

### Objectives

In this hands-on lab, you will learn how to:

- Download and load the stations\_info\_tx\_v31.0e.txt dataset (4.6 GB).
- Create Hive tables and perform queries using Hadoop Cluster.
- Conduct tempo-spatial analysis of maximum temperature patterns.
- Visualize regional trends in Excel.
- Understand the application of Big Data tools in climate research.

### Platform Specs

- IBM Bluemix BigInsights
- CPU Speed: 2.4 GHz
- # of CPU cores: 4 cores per node
- # of nodes: 3 nodes
- Total Memory Size: 24 GB (8 GB per node x 3 nodes)

## Step 1: Load Data to Hadoop Cluster

Log in to the remote cluster and download the data file:

```
$ ssh your_user@144.24.46.199
```

```
$ your_user@144.24.46.199 enter your password:
```

```
-bash-4.2$ wget -O stations_info_tx_v31.0e.txt https://knmi-ecad-  
assets-prd.s3.amazonaws.com/ensembles/data/stations_info_tx_v31.0e.txt
```

```
-bash-4.2$ wget -O stations_info_tx_v31.0e.txt https://knmi-ecad-assets-prd.s3.amazonaws.com/ensemb  
les/data/stations_info_tx_v31.0e.txt  
--2025-04-07 02:41:46-- https://knmi-ecad-assets-prd.s3.amazonaws.com/ensembles/data/stations_info  
_tx_v31.0e.txt  
Resolving knmi-ecad-assets-prd.s3.amazonaws.com (knmi-ecad-assets-prd.s3.amazonaws.com)... 52.218.6  
2.10, 52.218.40.26, 52.92.2.81, ...  
Connecting to knmi-ecad-assets-prd.s3.amazonaws.com (knmi-ecad-assets-prd.s3.amazonaws.com)|52.218.  
62.10|:443... connected.  
HTTP request sent, awaiting response... 200 OK  
Length: 433591 (423K) [text/plain]  
Saving to: 'stations_info_tx_v31.0e.txt'  
  
100%[=====>] 433,591      610KB/s   in 0.7s  
2025-04-07 02:41:47 (610 KB/s) - 'stations_info_tx_v31.0e.txt' saved [433591/433591]
```

Upload the file to the Stations directory of HDFS:

```
hdfs dfs -mkdir Stations
```

```
hdfs dfs -mkdir Stations/stations_info/
```

```
-bash-4.2$ hdfs dfs -mkdir Stations  
-bash-4.2$ hdfs dfs -mkdir Stations/stations_info/  
-bash-4.2$
```

```
hdfs dfs -put stations_info_tx_v31.0e.txt Stations/stations_info/
```

```
hdfs dfs -ls Stations/stations_info/
```

```
rm: cannot remove 'Stations': Is a directory  
-bash-4.2$ hdfs dfs -put stations_info_tx_v31.0e.txt Stations/stations_info/  
-bash-4.2$ hdfs dfs -ls Stations/stations_info/  
Found 1 items  
-rw-r--r-- 3 kalilai2 hdfs 433591 2025-04-29 02:56 Stations/stations_info/  
stations_info_tx_v31.0e.txt  
-bash-4.2$
```

## Step 2: Create Hive Table

Start Beeline CLI and switch to your assigned database.

```
$ beeline
```

```
use [your_database];
```

**If you do not have a database yet, create one:**

```
CREATE DATABASE your_database_name;
```

**Create the external Hive table by running these commands:**

```
DROP TABLE IF EXISTS max_temp_data;
```

```
CREATE EXTERNAL TABLE IF NOT EXISTS max_temp_data(station_id STRING,
```

```
    station_name STRING,
```

```
    country STRING,
```

```
    latitude DOUBLE,
```

```
    longitude DOUBLE,
```

```
    elevation DOUBLE,
```

```
    start_date STRING,
```

```
    end_date STRING)
```

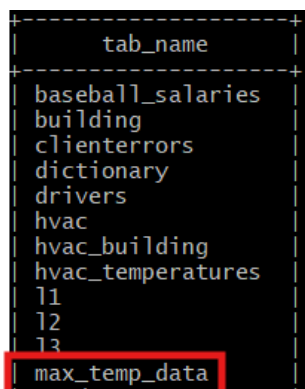
```
ROW FORMAT DELIMITED FIELDS TERMINATED BY '|'
```

```
STORED AS TEXTFILE LOCATION '/user/your_user/Stations/stations_info/'
```

```
TBLPROPERTIES ("skip.header.line.count"="1");
```

**After creating the table, check if the table is created and confirm the structure of the columns using these commands:**

```
SHOW tables;
```



tab_name
baseball_salaries
building
clienterrors
dictionary
drivers
hvac
hvac_building
hvac_temperatures
l1
l2
l3
max_temp_data

```
DESCRIBE max_temp_data;
```

col_name	data_type	comment
station_id	string	
station_name	string	
country	string	
latitude	double	
longitude	double	
elevation	double	
start_date	string	
end_date	string	

To make sure the data is loaded correctly, run:

```
SELECT * FROM max_temp_data LIMIT 5;
```

max_temp_data.station_id	max_temp_data.station_name	max_temp_data.country	max_temp_data.latitude	max_temp_data.longitude	max_temp_data.elevation	max_temp_data.start_date	max_temp_data.end_date
1	Vaexjoe	SWEDEN	56.87	14.8	166.0	1950-01-01	2006-12-31
2	Falun	SWEDEN	60.62	15.62	160.0	1950-01-01	2024-12-31
3	Stensele	SWEDEN	65.07	17.15	325.0	1950-01-01	2006-12-31
4	Linköping	SWEDEN	58.4	15.53	93.0	1950-01-01	2024-12-31
5	Linköping-Malmslätt	SWEDEN	58.4	15.53	93.0	1950-01-01	2024-12-31

### Step 3: Run Hive Queries

Run these queries for analysis:

Find out which countries had the highest average elevations of weather stations:

```
SELECT country, AVG(elevation) AS avg_elevation
FROM max_temp_data
GROUP BY country
ORDER BY avg_elevation DESC
LIMIT 10;
```

country	avg_elevation
ARMENIA	1955.075
KYRGYZSTAN	1398.5
IRAN, ISLAMIC REPUBLIC OF	1355.142857142857
BOSNIA AND HERZEGOVINA	1348.5
SWITZERLAND	1018.972972972973
TAJIKISTAN	934.0
NORTH MACEDONIA	779.5
AUSTRIA	758.0833333333334
SAUDI ARABIA	689.0
TÜRKIYE	568.6090909090908

10 rows selected (8.929 seconds)  
0: jdbc:hive2://bigdaiun0.sub03291929060.trai>

**Which countries have the most stations (more spatial coverage)?**

```
SELECT country, COUNT(*) AS num_stations
FROM max_temp_data
GROUP BY country
ORDER BY num_stations DESC
LIMIT 5;
```

country	num_stations
GERMANY	1013
RUSSIAN FEDERATION	318
SWEDEN	275
ITALY	204
NORWAY	183

5 rows selected (8.858 seconds)  
0: jdbc:hive2://bigdaiun0.sub03291929060.trai>

**What are the top 5 stations at the highest elevations?**

```
SELECT station_id, country, elevation
FROM max_temp_data
ORDER BY elevation DESC
LIMIT 5;
```

station_id	country	elevation
2073	ITALY	3480.0
2941	ARMENIA	3223.0
15	AUSTRIA	3109.0
58	GERMANY	2964.0
878	ITALY	2600.0

## Step 4: Downloading Data to Your Personal Computer

Run an **INSERT OVERWRITE DIRECTORY** query to save the table's content into a file in HDFS:

```
INSERT OVERWRITE DIRECTORY '/user/your_user/max_temp_export/'  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ','  
SELECT * FROM max_temp_data;
```

This will export all the data from your `max_temp_data` table into a new HDFS directory `/user/your_user/max_temp_export/`.

Open another terminal to download the output file at the HDFS path.

```
$ ssh your_user@144.24.46.199  
$ your_user@144.24.46.199 enter your password:  
xpthl@Kats_Laptop MINGW64 /  
$ ssh kalilai2@144.24.46.199  
kalilai2@144.24.46.199's password:  
Last login: Tue Apr 29 03:18:25 2025 from 172.56.235.234  
-bash-4.2$
```

Locate the file.

```
hdfs dfs -ls /user/your_user/max temp export/  
-bash-4.2$ hdfs dfs -ls /user/kalilai2/max_temp_export/  
Found 2 items  
drwxr-xr-x - kalilai2 hdfs 0 2025-04-29 03:45 /user/kalilai2/max_temp  
_export/.hive-staging_hive_2025-04-29_03-45-19_607_2775965156255915969-2211  
-rw-r--r-- 3 kalilai2 hdfs 379948 2025-04-30 08:02 /user/kalilai2/max_temp  
_export/000000_0  
-bash-4.2$
```

Export manually.

```
hdfs dfs -copyToLocal /user/your_user/max_temp_export/000000_0  
max_temp_data.csv
```

Then transfer to local machine:

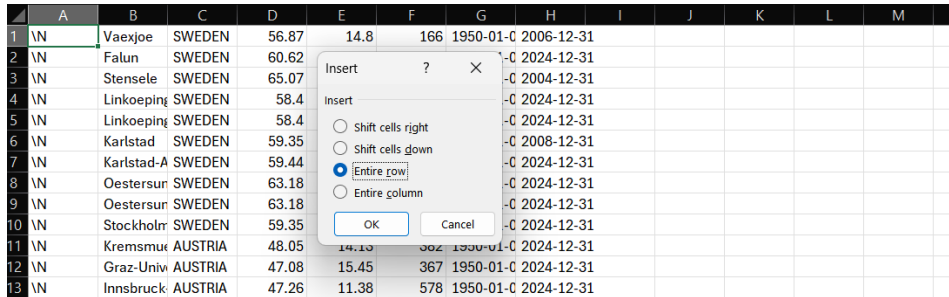
```
$ scp  
your_user@your_ip_address:/home/your_local_machine/max_temp_data.csv  
~/Downloads/
```

```
$ scp kalilai2@144.24.46.199:/home/kalilai2/max_temp_data.csv ~/Downloads/  
kalilai2@144.24.46.199's password:  
max_temp_data.csv 100% 371KB 1.0MB/s 00:00
```

## Step 5: Visualize Data

Open the downloaded file through Excel.

Right click on the first row and insert a new row above to create column headers.

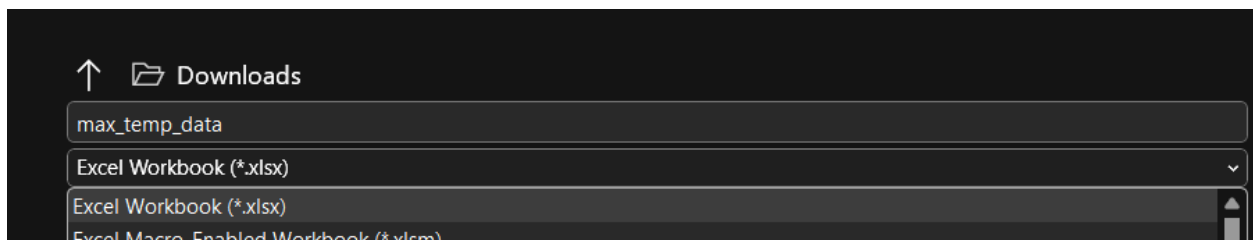


	A	B	C	D	E	F	G	H	I	J	K	L	M
1	VN	Vaexjoe	SWEDEN	56.87	14.8	166	1950-01-01	2006-12-31					
2	VN	Falun	SWEDEN	60.62				2024-12-31					
3	VN	Stensele	SWEDEN	65.07				2004-12-31					
4	VN	Linköping	SWEDEN	58.4				2024-12-31					
5	VN	Linköping	SWEDEN	58.4				2024-12-31					
6	VN	Karlstad	SWEDEN	59.35				2008-12-31					
7	VN	Karlstad-A	SWEDEN	59.44				2024-12-31					
8	VN	Oestersund	SWEDEN	63.18				2024-12-31					
9	VN	Oestersund	SWEDEN	63.18				2024-12-31					
10	VN	Stockholm	SWEDEN	59.35				2024-12-31					
11	VN	Kremsmühl	AUSTRIA	48.05	14.13	362	1950-01-01	2024-12-31					
12	VN	Graz-Universität	AUSTRIA	47.08	15.45	367	1950-01-01	2024-12-31					
13	VN	Innsbruck	AUSTRIA	47.26	11.38	578	1950-01-01	2024-12-31					

Insert column headers: station\_id, station\_name, country, latitude, longitude, elevation, start\_date, end\_date

	A	B	C	D	E	F	G	H
1	station_id	station_name	country	latitude	longitude	elevation	start_date	end_date
2	VN	Vaexjoe	SWEDEN	56.87	14.8	166	1950-01-01	2006-12-31
3	VN	Falun	SWEDEN	60.62	15.62	160	1950-01-01	2024-12-31

Go to File > Save As, and select file type as 'Excel Workbook (\*.xlsx)'.



Create two new columns called 'fixed\_start\_date' and 'fixed\_end\_date'. Use this formula, =DATEVALUE(cell), and select the cell containing the corresponding dates.

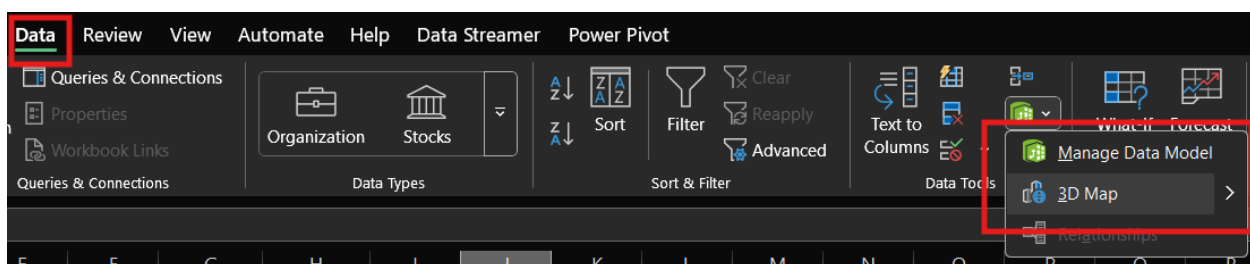
Drag the corner of the formulated cell to copy the formula to the rest of the column.

start_date	end_date	fixed_start_date	fixed_end_date
1950-01-01	2006-12-31	1/1/1950	=DATEVALUE(H2)
1950-01-01	2024-12-31	1/1/1950	12/31/2024
1950-01-01	2004-12-31	1/1/1950	12/31/2004

start_date	end_date	fixed_start_date	fixed_end_date
1950-01-01	2006-12-31	1/1/1950	12/31/2006
1950-01-01	2024-12-31	1/1/1950	12/31/2024
1950-01-01	2004-12-31	1/1/1950	12/31/2004
1950-01-01	2024-12-31	1/1/1950	12/31/2024
1950-01-01	2024-12-31	1/1/1950	12/31/2024

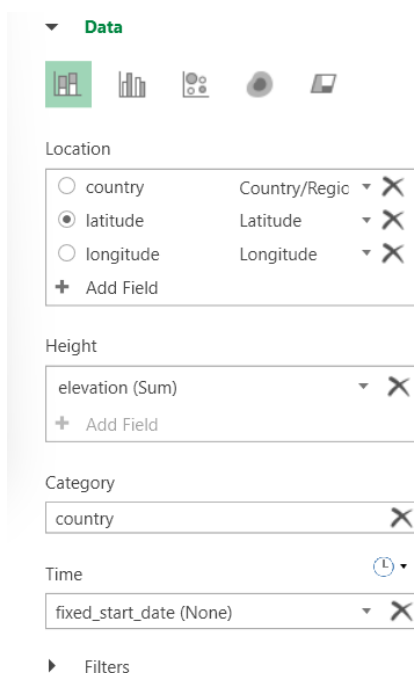
**NOTE:** You may also need to review the data to make sure everything is formatted correctly and fix any errors.

Next, go to the 'Data' tab, then 'Data Model' > 3D Map.



On the 3D Map screen, in the Data pane, select 'latitude' as the Location.

Drag 'elevation' to Height, 'country' to Category, and 'fixed\_start\_date' to Time.





It should generate this map with the ability to see changes over time.



## References

1. URL of Data Source,  
[https://surfobs.climate.copernicus.eu/dataaccess/access\\_eobs.php#datafiles](https://surfobs.climate.copernicus.eu/dataaccess/access_eobs.php#datafiles)
2. URL of your Github, <https://github.com/kalilai2/CIS-4560-01-Group-Project.git>