



CSCI 5408 – Data Management, Warehousing, Analytics

Assignment 5

Work done by,

Name: Guturu Rama Mohan Vishnu

Banner ID: B00871849

Email: rm286720@dal.ca

DECLARATION

I, Guturu Rama Mohan Vishnu, declare that in assignment 3 of CSCI 5408 course, I declare that all the work done was done by myself and I have not collaborated with anyone for the assignment.

Problem #1

Business Intelligence reporting using Cognos

Step 3 : Measurable Facts and Dimensions

Facts are numeric values that represent a specific business aspect. A fact is a metric that can be justified by a number of different factors. A foreign key relationship will exist between each fact table and the dimension table. Fact table is the center of the star schema. Dimensions are the qualifying characteristics that provide additional perspectives to a given fact. Dimension table holds the data in a way that can be used to analyze the data in fact table.

From the dataset I have downloaded from Kaggle, it is identified that the dataset is about the climate weather which covers hourly weather data from various weather stations in Brazil. So, the factor in the schema is the fact_table and the dimension tables are

- air relative humidity,
- air temperature,
- atmospheric pressure at station height,
- atmospheric pressure max in the previous hour,
- atmospheric pressure min in the previous hour,
- date.
- hour,
- dew point temperature,
- index,
- location,
- radiation,
- relative humidity max in the previous hour,

- relative humidity min in the previous hour,
- max temperature in the previous hour,
- min temperature in the previous hour,
- dew temperature max in the previous hour,
- dew temperature min in the previous hour,
- total precipitation,
- wind direction,
- wind rajada maxima and
- wind speed,

Step 4 : Cleaning and Formatting

Cleaning the dataset requires multiple steps to be processed and requires multiple ways of cleaning the data. Since the dataset is too large, I have combined all the 5 files and randomized the 50000 rows and extracted 5000 rows out of it and made a new excel sheet.

- To start with, there are many rows with integer type values in it and empty cells in those rows can be replaced by value '0' to maintain the consistency.

Before cleaning:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
2	2764	6165911	75395	30-04-2013		7:00	0	883.4	883.7	883.4		15	12.3	15.6	14.9	12.4	10.8	84	73	84
3	6314	6004621	19141	27-11-2013		9:00	0	914	914	913.4	56	19	18.2	19	18.7	18.2	17.9	95	95	95
4	3663	6092945	47434	19-02-2013		6:00	0	996.6	996.7	996.5		22.5	20.9	23.2	22.3	21.7	20.6	92	90	91
5	8421	976501	155876	21-11-2007		20:00	0	961.1	961.1	960.9	215	24.6	18.1	25.1	24.6	18.2	17.6	67	64	67
6	6771	13614390	7110	4/10/2020		3:00	0	1015.4	1015.7	1015.4		23.9	20.9	24.8	23.9	21.1	20.9	83	79	83
7	11166	6112139	66618	26-05-2013		12:00	0.2	972.8	972.8	972.4	1050	16.1	13.8	16.1	12.8	14	12	95	85	86
8	12159	12489579	39356	19-04-2019		21:00	0	1010	1010	1009.8	40	23.6	20.2	25.3	23.6	20.3	19.8	81	72	81
9	11915	7895514	3954	13-01-2015		0:00	0					23.1	18.9	24.9	23.1	19.1	18.6	77	69	77
10	8001	13237675	270398	13-11-2019		14:00	0	933.5	933.8	933.5	790	24.2	18.6	24.4	22.7	19.5	18.4	80	70	71
11	6171	14226628	160763	16-12-2020		9:00	0	915.5	915.5	915	33	21.3	14.8	21.4	20.8	15.8	14.7	73	66	66
12	12489	8995386	47218	16-04-2016		10:00	0	952.3	952.3	951.6	52	19.3	15.7	19.3	17.9	15.7	14.9	83	80	80
13	6744	7840391	413698	31-07-2014		2:00	0	912.9	912.9	912.8		11.4	11	11.8	11.3	11.6	11	99	98	98
14	10177	14085569	118236	17-06-2020		12:00	0.2	921.1	921.1	920.6	370	13.7	13.4	14.1	13.5	13.8	13.3	98	98	98
15	14773	6520047	230112	9/4/2013		0:00	0.2	831.5	831.5	830.9		13.3	13.2	13.5	13.3	13.4	13.2	100	99	99
16	10659	10981750	380899	22-08-2017		7:00	0	993.5	993.6	993.4		19.3	13.7	19.4	19.2	13.9	13.6	71	69	70
17	10764	3847310	339918	7/7/2010		6:00	0	919.1	919.6	919.1		15.7	7.5	16.6	15.1	7.5	7	60	54	58
18	5360	11447623	131823	15-10-2018		6:00	0	1010.3	1010.4	1010		20	17.7	20.1	19.9	17.9	17.7	87	87	87
19	2560	13379796	322517	26-10-2019		5:00	0	934	934.7	934		20	18.2	20.1	19.7	18.3	17.9	90	89	89
20	13053	11418528	117709	29-10-2018		16:00	0	928.3	929.2	928.3	2886	28.1	18.8	28.1	25.4	19.4	18	68	57	57
21	12255	9255796	142634	9/6/2016		23:00	0	955	955	954.5		10.1	8.8	11.5	10	9.7	8.5	91	87	91
22	908	9167796	114611	23-10-2016		8:00	0	964.3	964.4	963.9		18.4	14.9	18.7	18.4	14.9	14.9	80	79	80
23	2042	14152044	139689	7/11/2020		6:00	0	985.8	985.8	985.4		22.5	19.9	22.7	22.5	20.2	19.9	86	85	85
24	5528	15065610	95568	23-01-2021		0:00		922	922	921.5		21	16.2	21.9	20.7	16.4	16.2	76	71	74

After cleaning:

final_new_southeast - Excel

Sanjuna Konda

File Home Insert Page Layout Formulas Data Review View Help Tell me what you want to do

Clipboard Font Alignment Number Styles Cells Editing

Calibri 11 A A

B I U

General

Conditional Formatting Format as Table Cell Styles Insert Delete Format

AutoSum Fill Clear Sort & Find & Filter Select

T13 205

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
13	6744	7840391	413698	31-07-2014	2:00	0	912.9	912.9	912.8	0	11.4	11	11.8	11.3	11.6	11	99	98	98	
14	10177	14085569	118236	17-06-2020	12:00	0.2	921.1	921.1	920.6	370	13.7	13.4	14.1	13.5	13.8	13.3	98	98	98	
15	14773	6520047	230112	9/4/2013	0:00	0.2	831.5	831.5	830.9	0	13.3	13.2	13.5	13.3	13.4	13.2	100	99	99	
16	10659	10981750	380899	22-08-2017	7:00	0	993.5	993.6	993.4	0	19.3	13.7	19.4	19.2	13.9	13.6	71	69	70	
17	10764	3847310	339918	7/7/2010	6:00	0	919.1	919.6	919.1	0	15.7	7.5	16.6	15.1	7.5	7	60	54	58	
18	5360	11447623	131823	15-10-2018	6:00	0	1010.3	1010.4	1010	0	20	17.7	20.1	19.9	17.9	17.7	87	87	87	
19	2560	13379796	322517	26-10-2019	5:00	0	934	934.7	934	0	20	18.2	20.1	19.7	18.3	17.9	90	89	89	
20	13053	11418528	117709	29-10-2018	16:00	0	928.3	929.2	928.3	2886	28.1	18.8	28.1	25.4	19.4	18	68	57	57	
21	12255	9255796	142634	9/6/2016	23:00	0	955	955	954.5	0	10.1	8.8	11.5	10	9.7	8.5	91	87	91	
22	908	9167796	114611	23-10-2016	8:00	0	964.3	964.4	963.9	0	18.4	14.9	18.7	18.4	14.9	14.9	80	79	80	
23	2042	14152044	139689	7/11/2020	6:00	0	985.8	985.8	985.4	0	22.5	19.9	22.7	22.5	20.2	19.9	86	85	85	
24	5528	15065610	95568	23-01-2021	0:00	0	922	922	921.5	0	21	16.2	21.9	20.7	16.4	16.2	76	71	74	
25	4546	11556443	165625	28-11-2018	1:00	0	993.2	993.2	993.1	0	20.6	18.8	21	20.4	18.9	18.5	90	87	89	
26	4795	6672440	307500	2/2/2013	12:00	0	951.7	951.7	951.1	1407	24.7	19.8	24.8	23.1	20	19.3	81	74	74	
27	4372	11546467	155675	9/9/2018	14:00	0	904.6	904.9	904.6	2916	23.1	9.5	24	20.9	10.5	7.8	49	40	42	
28	14149	2573861	162614	25-07-2009	14:00	0.2	1016.5	1017.5	1016.5	0	20.5	18.2	20.5	19.6	18.2	17.6	89	86	87	
29	9017	12408777	18585	5/12/2019	10:00	0	1008.1	1008.1	1007.8	620	27.7	20.8	27.7	25.7	20.9	20.3	74	66	66	
30	5604	5628921	274641	4/1/2012	14:00	0	983	983.4	982.9	1603	28.2	20.6	28.3	27.4	21.6	19.5	69	60	64	
31	7376	153300	104724	18-01-2003	12:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
32	2981	1567417	105142	10/4/2008	22:00	0.2	906.1	906.2	905.8	0	19.8	18.1	21.1	19.8	18.5	17.3	90	80	90	
33	1490	13848008	60699	10/11/2020	0:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
34	14979	6828692	388757	14-05-2013	5:00	0	948.3	948.6	948.3	0	14.3	12.2	14.5	12.3	12.4	10.8	93	86	87	
35	9745	8589593	297414	14-07-2015	12:00	0	873	873.1	872.8	554	15.5	14.4	15.5	14.9	14.4	13.7	94	91	93	

final_new_southeast

Ready

Type here to search

11:49 PM 3/26/2022

- Next, there are few miscellaneous symbols in the data sheet like (Â). I am removing these symbols in the process of cleaning.

Before cleaning:

final_new_southeast - Excel

Sanjuna Konda

File Home Insert Page Layout Formulas Data Review View Help Tell me what you want to do

Clipboard Font Alignment Number Styles Cells Editing

Calibri 11 A A

B I U

General

Conditional Formatting Format as Table Cell Styles Insert Delete Format

AutoSum Fill Clear Sort & Find & Filter Select

A1 0

	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	hour	total precipitation (mm)	atmospheric pressure at station height (mb)	atmospheric pressure max. in the previous hour (mb)	atmospheric pressure min. in the previous hour (mb)	radiation (kJ/m2)	air temperature - dry bulb (Â°C)	dew point temperature (Â°C)	max. temperature in the previous hour (Â°C)	min. temperature in the previous hour (Â°C)	dew temperature max. in the previous hour (Â°C)	dew temperature min. in the previous hour (Â°C)	relative humidity max. in the previous hour (%)	relative humidity min. in the previous hour (%)	air relative humidity (Â° (gr))	wind direction (m/s)	wind speed (m/s)	station code		
2	30-04-2013	7:00	0	883.4	883.7	883.4	0	15	12.3	15.6	14.9	12.4	10.8	84	73	84	100	6.4	A555	
3	27-11-2013	9:00	0	914	914	913.4	56	19	18.2	19	18.7	18.2	17.9	95	95	95	67	3.6	A551	
4	19-02-2013	6:00	0	996.6	996.7	996.5	0	22.5	20.9	23.2	22.3	21.7	20.6	92	90	91	247	1.7	A623	
5	21-11-2007	20:00	0	961.1	961.1	960.9	215	24.6	18.1	25.1	24.6	18.2	17.6	67	64	67	49	4.1	A527	
6	4/10/2020	3:00	0	1015.4	1015.7	1015.4	0	23.9	20.9	24.8	23.9	21.1	20.9	83	79	83	191	5.8	A612	
7	26-05-2013	12:00	0.2	972.8	972.8	972.4	1050	16.1	13.8	16.1	12.8	14	12	95	85	86	120	4	A718	
8	19-04-2019	21:00	0	1010	1010	1009.8	40	23.6	20.2	25.3	23.6	20.3	19.8	81	72	81	192	5.4	A601	
9	13-01-2015	0:00	0	0	0	0	0	23.1	18.9	24.9	23.1	19.1	18.6	77	69	77	249	3.5	A557	
10	13-11-2019	14:00	0	933.5	933.8	933.5	790	24.2	18.6	24.4	22.7	19.5	18.4	80	70	71	236	4.5	A536	
11	16-12-2020	9:00	0	915.5	915.5	915	33	21.3	14.8	21.4	20.8	15.8	14.7	73	66	66	136	1.4	A553	
12	16-04-2016	10:00	0	952.3	952.3	951.6	52	19.3	15.7	19.3	17.9	15.7	14.9	83	80	80	36	1.6	A519	
13	31-07-2014	2:00	0	912.9	912.9	912.8	0	11.4	11	11.8	11.3	11.6	11	99	98	98	205	4.3	A613	
14	17-06-2020	12:00	0.2	921.1	921.1	920.6	370	13.7	13.4	14.1	13.5	13.8	13.3	98	98	98	65	5.7	A518	
15	9/4/2013	0:00	0.2	831.5	831.5	830.9	0	13.3	13.2	13.5	13.3	13.4	13.2	100	99	99	0	0	A610	
16	22-08-2017	7:00	0	993.5	993.6	993.4	0	19.3	13.7	19.4	19.2	13.9	13.6	71	69	70	126	3.6	A540	

final_new_southeast

Ready

Type here to search

1:33 AM 3/27/2022

After cleaning:

- The date column has values with different formats which might be an issue later, in case of querying or retrieving the data. So, I have changed that column to a single format of all values to maintain consistency.

Before cleaning:

final_new_southeast - Excel

Sanjuna Konda

File Home Insert Page Layout Formulas Data Review View Help Tell me what you want to do

Clipboard Font Alignment Number Styles Cells Editing

Calibri 11 A A Wrap Text General Conditional Formatting Format as Table Cell Styles Insert Delete Format AutoSum Fill Sort & Find & Filter Select

hour

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y
16	10659	1.1E+07	380899	22-08-2017	7.00	0	993.5	993.6	993.4	0	19.3	13.7	19.4	19.2	13.9	13.6	71	69	70	126	3.6	A540			
17	10764	3847310	339918	7/7/2010	6.00	0	919.1	919.6	919.1	0	15.7	7.5	16.6	15.1	7.5	7	60	54	58	89	5.8	A548			
18	5360	1.1E+07	318123	15-10-2018	6.00	0	1010.3	1010.4	1010	0	20	17.7	20.1	19.9	17.9	17.7	87	87	87	116	1	A603			
19	2560	1.3E+07	322517	26-10-2019	5.00	0	934	934.7	934	0	20	18.2	20.1	19.7	18.3	17.9	90	89	89	295	0.9	A633			
20	13053	1.1E+07	117709	29-10-2018	16.00	0	928.3	929.2	928.3	2886	28.1	18.8	28.1	25.4	19.4	18	68	57	57	99	7.4	A549			
21	12255	9255796	142634	9/6/2016	23.00	0	955	955	954.5	0	10.1	8.8	11.5	10	9.7	8.5	91	87	91	226	1.5	A728			
22	908	9167796	114611	23-10-2016	8.00	0	964.3	964.4	963.9	0	18.4	14.9	18.7	18.4	14.9	14.9	80	79	80	111	10.8	A707			
23	2042	1.4E+07	139689	7/11/2020	6.00	0	985.8	985.8	985.4	0	22.5	19.9	22.7	22.5	20.2	19.9	86	85	85	40	4.1	A540			
24	5528	1.5E+07	95568	23-01-2021	0.00	0	922	922	921.5	0	21	16.2	21.9	20.7	16.4	16.2	76	71	74	96	5.5	A533			
25	4546	1.2E+07	165625	28-11-2018	1.00	0	993.2	993.2	993.1	0	20.6	18.8	21	20.4	18.9	18.5	90	87	89	131	3.9	A522			
26	4795	6672440	307500	2/2/2013	12.00	0	951.7	951.7	951.1	1407	24.7	19.8	24.8	23.1	20	19.3	81	74	74	25	3.8	A519			
27	4372	1.2E+07	155675	9/9/2018	14.00	0	904.6	904.9	904.6	2916	23.1	9.5	24	20.9	10.5	7.8	49	40	42	11	6	A529			
28	14149	2573861	162614	25-07-2009	14.00	0.2	1016.5	1017.5	1016.5	0	20.5	18.2	20.5	19.6	18.2	17.6	89	86	87	30	6.6	A652			
29	9017	1.2E+07	18585	5/12/2019	10.00	0	1008.1	1008.1	1007.8	620	27.7	20.8	27.7	25.7	20.9	20.3	74	66	66	9	8.8	A602			
30	5604	5628921	274641	4/1/2012	14.00	0	983	983.4	982.9	1603	28.2	20.6	28.3	27.4	21.6	19.5	69	60	64	315	4.4	A540			
31	7376	153300	104724	18-01-2003	12.00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	A507
32	2981	1567417	105142	10/4/2008	22.00	0.2	906.1	906.2	905.8	0	19.8	18.1	21.1	19.8	18.5	17.3	90	80	90	131	3.3	A515			
33	1490	1.4E+07	60699	10/11/2020	0.00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	A541
34	14979	6828692	388757	14-05-2013	5.00	0	948.3	948.6	948.3	0	14.3	12.2	14.5	12.3	12.4	10.8	93	86	87	200	2.4	A713			
35	9745	8589593	297414	14-07-2015	12.00	0	873	873.1	872.8	554	15.5	14.4	15.5	14.9	14.4	13.7	94	91	93	44	4.2	A537			
36	6406	254577	92937	18-07-2004	9.00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	A706
37	12059	2552170	140923	5/4/2009	9.00	0	966.7	966.7	966.2	0	19.9	19.2	20.8	19.7	20	18.8	96	94	96	191	0	A718			
38	14346	9128476	86914	3/12/2015	10.00	0	906.1	906.3	905.8	529	21.7	19.5	21.7	18.6	19.8	18.2	98	87	87	133	2.8	A523			
39	13876	5261531	100298	29-02-2012	7.00	0	903.1	903.2	903.1	0	22.7	14.9	23	21.9	15	14.4	65	59	61	90	1.9	A708			
40	2704	1.5E+07	134043	7/3/2021	3.00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	A741
41	13772	4713644	288946	29-03-2011	7.00	0	919	919.3	919	0	16.6	15.9	17.1	16.2	16.1	15.3	96	92	96	306	1	A533			
42	3379	1.4E+07	111218	10/8/2020	23.00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	A766
43	14842	3789116	311734	18-04-2010	22.00	0	952.9	952.9	952.7	0	23.9	13.7	26.2	23.9	13.7	12.5	53	44	53	164	3.3	A726			
44	14569	4538617	209497	15-07-2011	6.00	0	950.8	951.3	950.8	0	19.5	12	20.8	19.2	12.4	11.2	63	57	62	132	5.4	A543			

Ready Average: 41757.28098 Count: 6001 Sum: 98380154 80%

Type here to search

After cleaning:

southeast.csv - Excel

Rama Mohan Vishnu Guturu

File Home Insert Page Layout Formulas Data Review View Help Tell me what you want to do

Clipboard Font Alignment Number Styles Cells Editing

Calibri 11 A A Date Conditional Formatting Format as Table Cell Styles Insert Delete Format AutoSum Fill Sort & Find & Filter Select

12/14/2000

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
292	5319	12/14/2000	15:00	0	1003	1003.8	1003	3655	33.3	19.6	33.3	32.1	21.1	19.4	50	44	44	15	8.2	3.4
293	882	6/12/2000	18:00	0	1010.9	1011	1010.7	1458	26	17.9	31.3	26	18.6	10.8	63	28	61	212	9.5	5.4
294	2385	12/7/2001	9:00	0	939.5	939.5	939	8	21.1	20.2	21.1	21	20.3	20.2	95	94	94	346	5	1.5
295	992	10/10/2001	8:00	0	941.1	941.1	940.4	0	20.5	19.1	20.8	20.5	19.1	19	92	90	92	75	3.5	2.3
296	4956	11/29/2000	12:00	0	1008.1	1008.2	1008	2557	30.9	20.3	30.9	29.7	20.3	19.2	55	51	53	6	10.3	5
297	2889	9/4/2000	9:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
298	549	5/29/2000	21:00	0	1012.2	1012.2	1011.5	6	18	13.1	19.7	18	13.5	13.1	73	67	73	198	5	2.3
299	2362	8/13/2000	10:00	0	1021.8	1021.8	1021.3	102	13.9	12.9	13.9	13.2	13	12.4	96	94	94	344	3.1	0.2
300	4368	11/5/2000	0:00	0	1006.7	1007	1006.6	0	23	20.9	23	22.9	21	20.9	89	88	88	197	5.3	2
301	2519	12/12/2001	23:00	0	939.3	939.4	938.7	0	22.4	19.8	22.6	22.2	20.3	19.7	89	85	85	96	2.1	0.9
302	1367	7/2/2000	23:00	0	1012	1012	1011.8	0	20.8	15.2	21.7	20.8	15.2	14.6	70	64	70	3	2.7	1.7
303	748	9/30/2001	4:00	0	944.9	946	944.9	0	18.6	15.4	19.2	18.6	15.5	15.4	82	79	81	119	8.6	3.5
304	3273	9/20/2000	9:00	0	1015.4	1015.4	1014.9	5	17.6	17	18.5	17.6	17.9	17	97	96	96	278	1	0.1
305	774	6/8/2000	6:00	0	1015.8	1016.1	1015.8	0	21.4	15.5	21.5	21	15.5	15.2	70	68	69	324	6.1	2.7
306	555	9/22/2001	3:00	0	944.7	945.2	944.7	0	18.8	14.8	18.9	18.7	15.3	14.8	81	78	78	134	4.9	1.8
307	3629	1/28/2001	5:00	0	1007	1007.2	1006.9	0	25.2	22.7	25.2	24.3	23	22.7	93	86	86	333	4.1	1.7
308	107	5/11/2000	11:00	0	1018.1	1018.1	1017.6	595	21.5	18.1	21.6	18.9	18.6	17.6	93	81	81	89	0	0.4
309	2747	8/29/2000	11:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
310	390	9/15/2001	6:00	0	937	937.6	936.8	0	18.8	18.4	18.8	18.4	18.4	17.9	97	97	97	347	5.8	2.1
311	3533	10/1/2000	5:00	0	1014.4	1015.2	1014.4	0	20.2	18.8	20.5	20.2	18.9	18.7	92	90	92	56	0.4	0
312	2436	12/9/2001	12:00	0	941.7	941.7	941.5	1521	24.9	20.1	25.2	24	20.5	19.5	78	73	75	52	5.2	2.5
313	733	6/6/2000	13:00	0	1020.7	1021	1020.7	1817	24.2	14.3	24.2	22.8	14.5	14.1	59	54	54	345	8.4	3.9
314	1630	7/13/2000	22:00	0	1020.9	1020.9	1020.6	0	17.9	13.6	18.1	17.9	13.6	13	76	72	76	188	3.1	1.8

Ready Average: 1/26/2001 Count: 3000 100%

Type here to search

- There were many values of -9999 in the dataset. So, I removed them and replaced them with 0.

Before cleaning:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
4	2	5/7/2000	2:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
5	3	5/7/2000	3:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
6	4	5/7/2000	4:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
7	5	5/7/2000	5:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
8	6	5/7/2000	6:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
9	7	5/7/2000	7:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
10	8	5/7/2000	8:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
11	9	5/7/2000	9:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
12	10	5/7/2000	10:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
13	11	5/7/2000	11:00	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	-9999	SE
14	12	5/7/2000	12:00	0	1012.5	1012.5	1012	219	22	18	22	20.8	18.9	17.8	87	77	78	30	1.4	0.7	SE
15	13	5/7/2000	13:00	0	1013	1013	1012.6	297	22.4	18	22.5	22	18	17.2	77	74	76	212	2.2	0.5	SE
16	14	5/7/2000	14:00	0	1012.7	1013	1012.7	527	22.6	17.6	22.8	22.4	18.3	17.6	77	73	73	213	4.6	2.5	SE
17	15	5/7/2000	15:00	0	1012.2	1012.7	1012.2	501	22.6	18	22.8	22.5	18	17.2	75	71	75	223	5.8	2.9	SE
18	16	5/7/2000	16:00	0	1012	1012.3	1012	549	22.1	17.4	22.6	21.8	18.1	17.4	77	74	74	239	8.7	3.7	SE
19	17	5/7/2000	17:00	0	1011.7	1012	1011.6	506	22.5	17	22.6	22.1	17.4	16.6	74	70	71	216	6.5	3.4	SE
20	18	5/7/2000	18:00	0	1011.8	1011.8	1011.6	412	22.9	17.3	22.9	22.5	17.3	16.7	71	69	71	200	5.2	2.8	SE
21	19	5/7/2000	19:00	0	1012.2	1012.2	1011.8	332	22.5	17.2	23	22.5	17.4	16.6	72	67	72	198	4.8	3.2	SE
22	20	5/7/2000	20:00	0	1012.5	1012.5	1012.1	142	22	16.9	22.6	22	17.6	16.8	74	71	73	214	4.6	1.5	SE
23	21	5/7/2000	21:00	0	1013.1	1013.1	1012.6	1	21.7	17.4	22	21.7	17.5	16.7	77	72	77	155	3.1	0.6	SE
24	22	5/7/2000	22:00	0	1013.7	1013.8	1013	-9999	21.7	17.6	21.9	21.6	17.7	17.4	78	76	77	199	3.3	1.6	SE

After cleaning:

The screenshot shows an Excel spreadsheet with the following data (rows 4 to 26):

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
4	2	5/7/2000	2:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
5	3	5/7/2000	3:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
6	4	5/7/2000	4:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
7	5	5/7/2000	5:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
8	6	5/7/2000	6:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
9	7	5/7/2000	7:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
10	8	5/7/2000	8:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
11	9	5/7/2000	9:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
12	10	5/7/2000	10:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
13	11	5/7/2000	11:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	SE
14	12	5/7/2000	12:00	0	1012.5	1012.5	1012	219	22	18	22	20.8	18.9	17.8	87	77	78	30	1.4	0.7	SE
15	13	5/7/2000	13:00	0	1013	1013	1012.6	297	22.4	18	22.5	22	18	17.2	77	74	76	212	2.2	0.5	SE
16	14	5/7/2000	14:00	0	1012.7	1013	1012.7	527	22.6	17.6	22.8	22.4	18.3	17.6	77	73	73	213	4.6	2.5	SE
17	15	5/7/2000	15:00	0	1012.2	1012.7	1012.2	501	22.6	18	22.8	22.5	18	17.2	75	71	75	223	5.8	2.9	SE
18	16	5/7/2000	16:00	0	1012	1012.3	1012	549	22.1	17.4	22.6	21.8	18.1	17.4	77	74	74	239	8.7	3.7	SE
19	17	5/7/2000	17:00	0	1011.7	1012	1011.6	506	22.5	17	22.6	22.1	17.4	16.6	74	70	71	216	6.5	3.4	SE
20	18	5/7/2000	18:00	0	1011.8	1011.8	1011.6	412	22.9	17.3	22.9	22.5	17.3	16.7	71	69	71	200	5.2	2.8	SE
21	19	5/7/2000	19:00	0	1012.2	1012.2	1011.8	332	22.5	17.2	23	22.5	17.4	16.6	72	67	72	198	4.8	3.2	SE
22	20	5/7/2000	20:00	0	1012.5	1012.5	1012.1	142	22	16.9	22.6	22	17.6	16.8	74	71	73	214	4.6	1.5	SE
23	21	5/7/2000	21:00	0	1013.1	1013.1	1012.6	1	21.7	17.4	22	21.7	17.5	16.7	77	72	77	155	3.1	0.6	SE
24	22	5/7/2000	22:00	0	1013.7	1013.8	1013	0	21.7	17.6	21.9	21.6	17.7	17.4	78	76	77	199	3.3	1.6	SE
25	23	5/7/2000	23:00	0	1014.3	1014.4	1013.7	0	21.8	16.9	22.2	21.7	17.6	16.4	77	70	74	224	6.8	1.1	SE
26	24	5/8/2000	0:00	0	1014.4	1014.5	1014.2	0	21.5	17.7	21.8	21.4	17.7	16.8	79	74	79	127	3.7	1.9	SE

- After the above step, if a row has all 0 values in it, then I am deleting the row from the dataset.
- There are no duplicate values in the rows I extracted from the original dataset. So, there is no need to delete the duplicate rows.

After performing all these steps, the dataset is cleaned and it doesn't have any unnecessary or unwanted data in it.

Step 5 : Create Cognos account

I have created the IBM Cognos account and have imported all the dimension tables along with the fact table.

The screenshot displays the IBM Cognos Analytics web application. The browser's address bar shows the URL `us1.ca.analytics.ibm.com/bi/?perspective=content`. The application header includes the text "IBM Cognos Analytics with Watson" and a search bar. The main content area is titled "Content" and features three tabs: "My content", "Team content", and "Samples". The "My content" tab is active, showing a grid of imported CSV files. Each file card displays the file name, a "Last Accessed" timestamp, and a "CSV" icon with an upload arrow. The files shown are:

- a5 (Last Accessed: 3/29/2022, 12:39 AM)
- air relative humidity.csv (Last Accessed: 3/28/2022, 9:18 PM)
- air temperature.csv (Last Accessed: 3/28/2022, 9:18 PM)
- atmospheric pressure.csv
- date and hour.csv
- dew point temperature.csv
- fact_table.csv

A user profile dropdown menu is open on the right side of the screen, showing the user's name "Rama Mohan Vishnu Gut...", email "rm286720@dal.ca", and a "Log out" button. The Windows taskbar at the bottom shows the system clock as 2:40 AM on 3/29/2022.

Step 6 : Create Star Schema/Snowflake Schema

I have created a star schema based on my understanding of the domain. As already said, there are multiple dimensions for the fact table and all these dimensions have 1 : N relationship with the fact table.

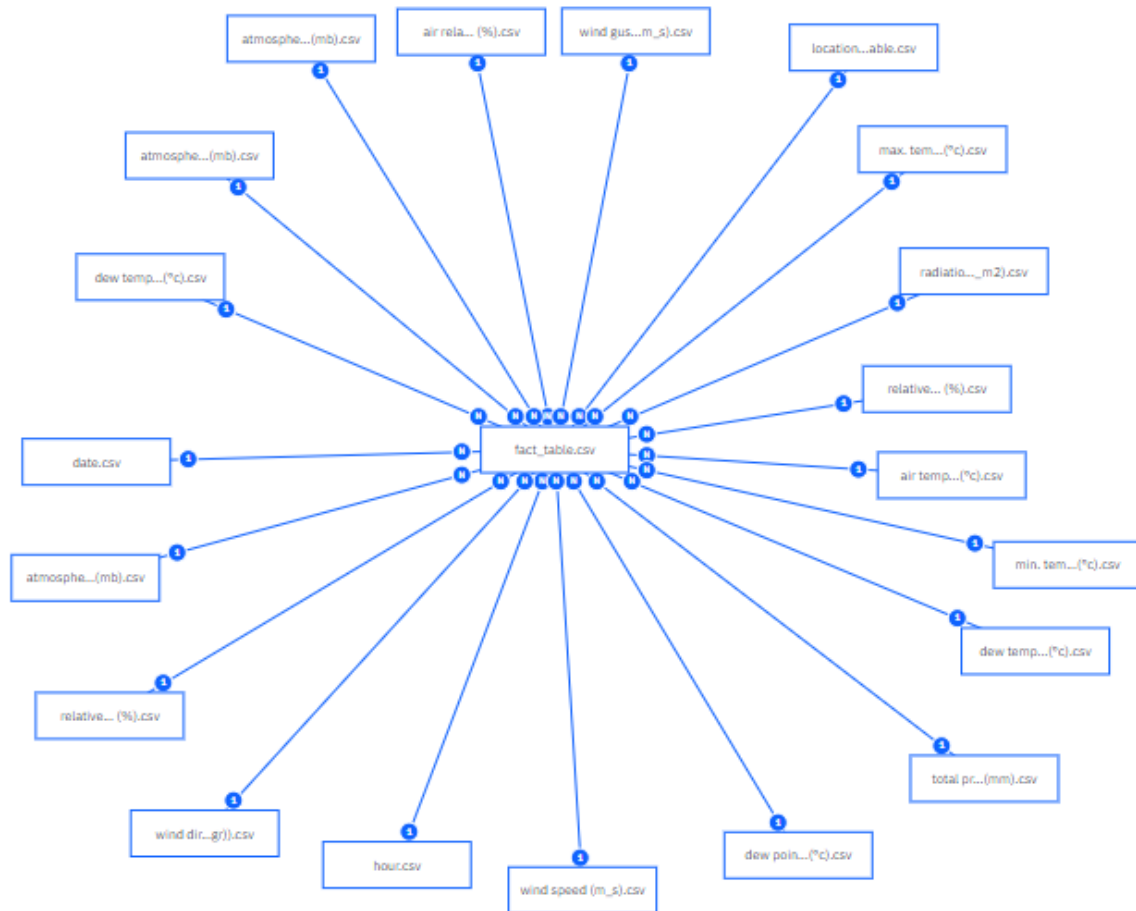
- ✓ Location – A region's climate is inextricably tied to the weather that occurs there. If a location experiences very cold weather for the majority of the year, it is considered to have a cold climate since the average temperature over a lengthy period of time is primarily cold.
- ✓ Wind – Wind transports moisture, as well as hot and cold air, into the atmosphere, influencing weather patterns. As a result, a change in wind causes a change in weather.
- ✓ Temperature – This is the main factor of all the factors for any weather. To be brief, temperature causes climate and climate is the result of temperature. So, they are bonded very closely.
- ✓ Relative Humidity – If the temperature increases, the humidity increases and climate gets hot. If the temperature decreases, humidity decreases and climate gets cold.
- ✓ Total Precipitation – Climate change has the potential to alter precipitation intensity and frequency. Warmer oceans cause more water to evaporate into the atmosphere. Heavy rain and snowstorms, for example.
- ✓ Atmospheric Pressure – Atmospheric pressure is a weather indicator. Cloudiness, wind, and precipitation are common when a low-pressure system moves into a region. Fair and quiet weather is frequently associated with high-pressure systems.
- ✓ Date and Hour – Weather changes every single day due to winds and storms. Weather also changes with seasons. In summer season, the weather is hot; in

winter, the weather is too cold and in rainy season, it is just below the room temperature.

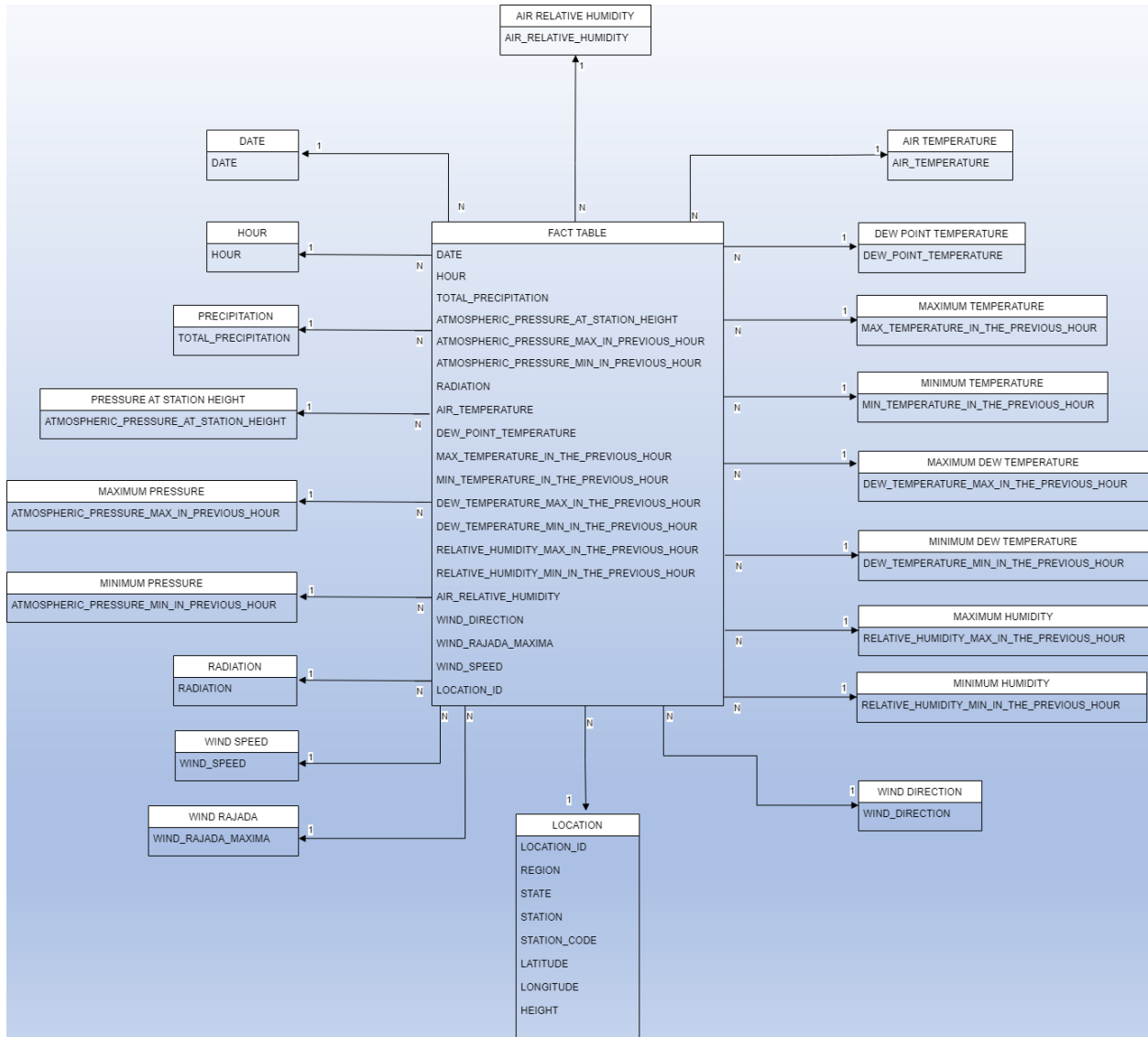
- ✓ Radiation – Radiation is very important to weather. The amount of radiation defines the temperature of that area of that city and the temperature directly affects the climate at that particular time.

Step 7:

The schema I created using IBM Cognos is provided below.



I have also built the dimensional modelling for this dataset using draw.io. Below is the model I built.

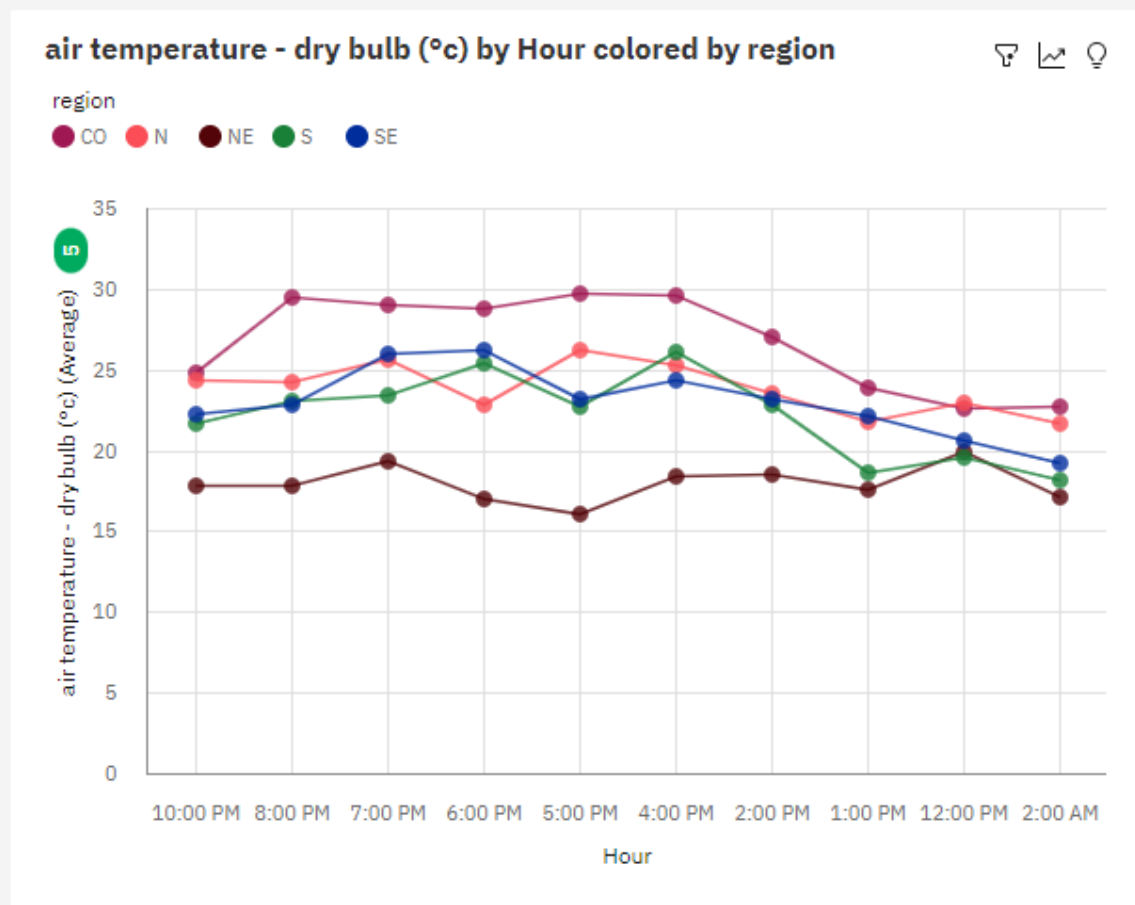


Step 8 : Visual Analysis of the data

In order for anyone to understand the dataset easily and in a better way, I created few visualizations which might help to look and understand better, rather than looking at thousands of rows of data in excel sheets.

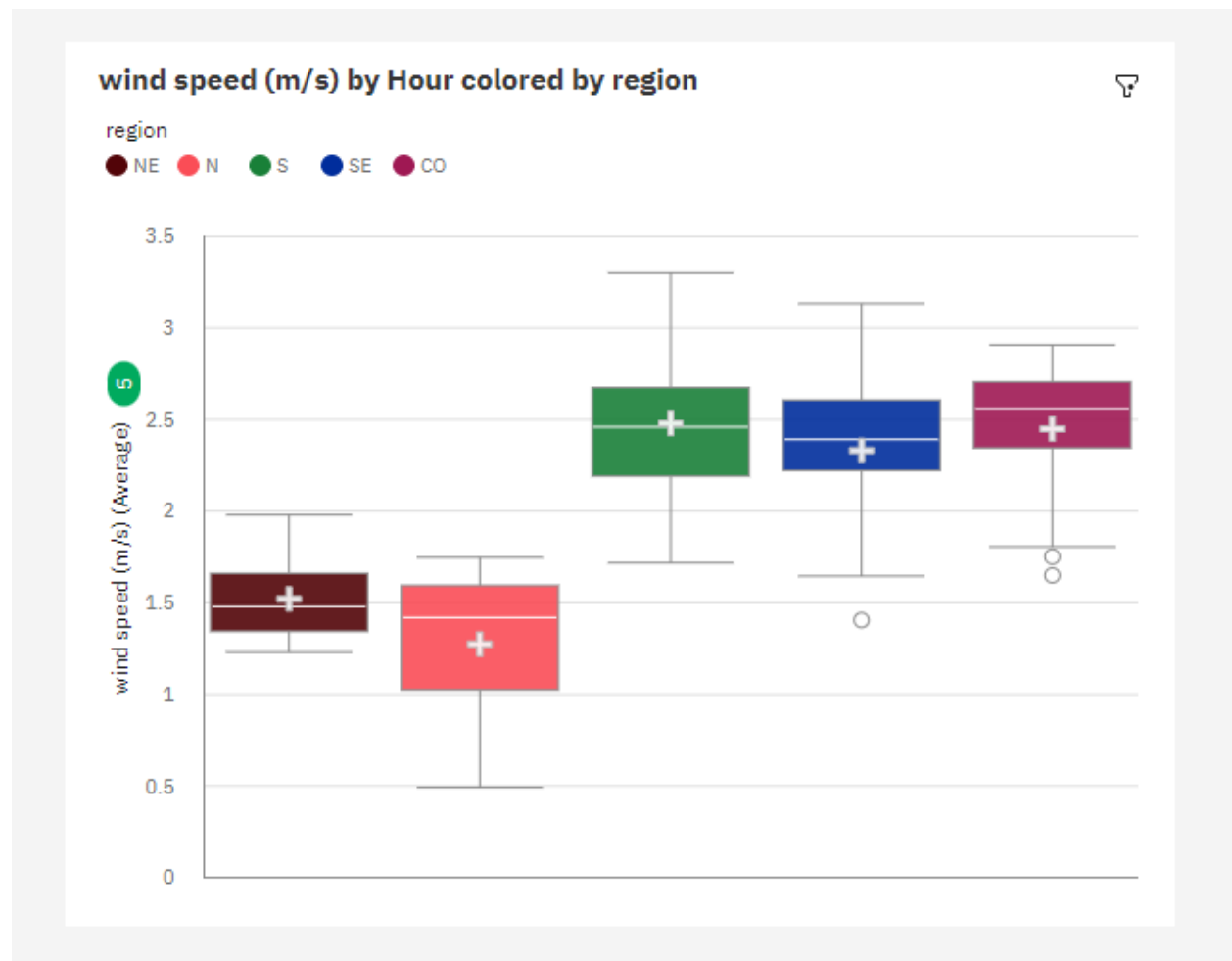
Visual 1 – Air temperature by Hour colored by Region

This visual represents the air temperature of all the regions based on hour by intervals. From this visual, we can understand that the value of air temperature is most unusual when the values of region are NE and CO. And also, the values of air temperature are most unusual when the values of hour are 02:00, 13:00 and 19:00.



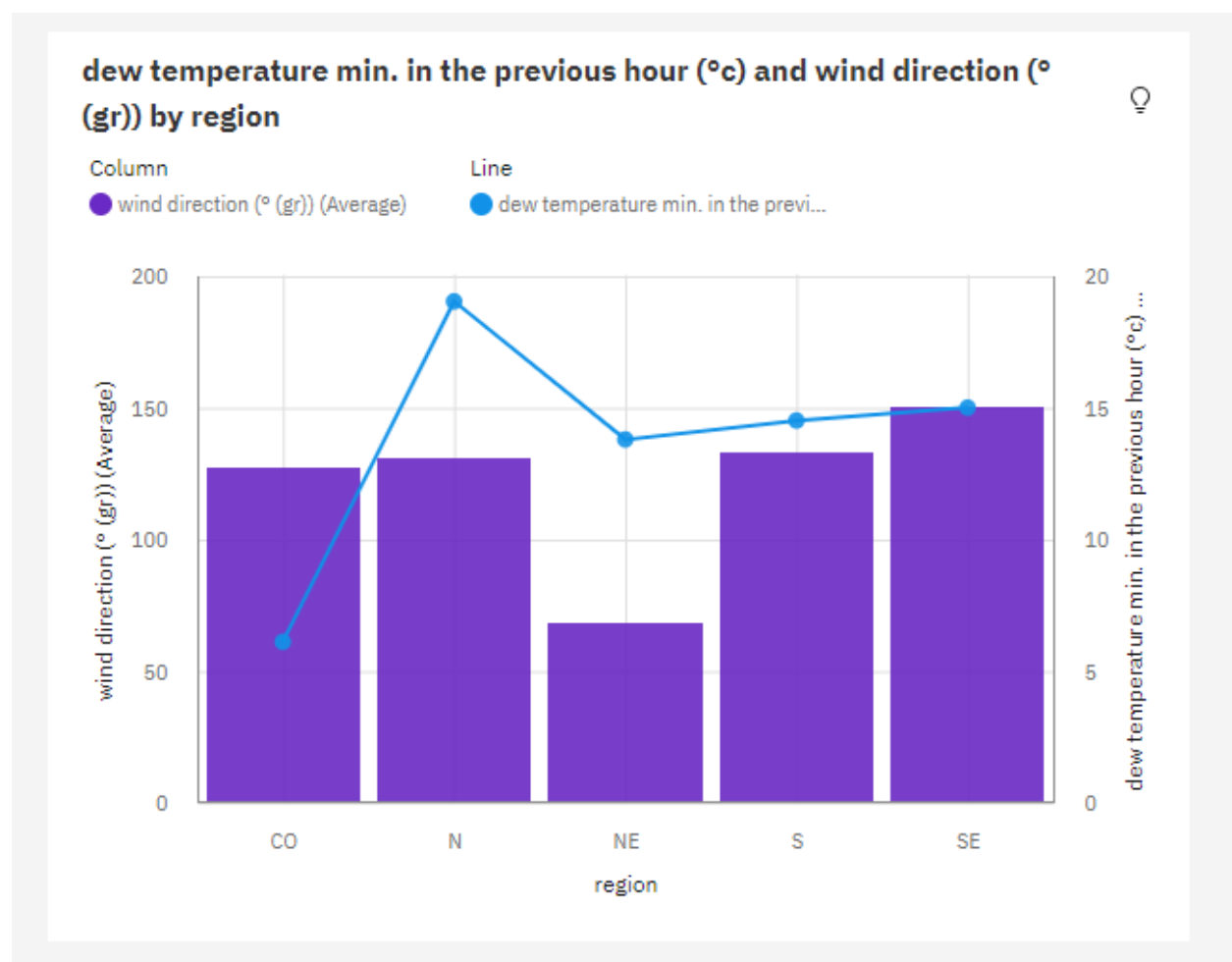
Visual 2 – Wind speed by Hour colored by Region

This visual represents the relative wind speeds separated by hours and differentiated by regions. As we can see, the average wind speed is higher in CO region and lower in N region. CO, SE and S have almost the same average wind speed with minor differences.



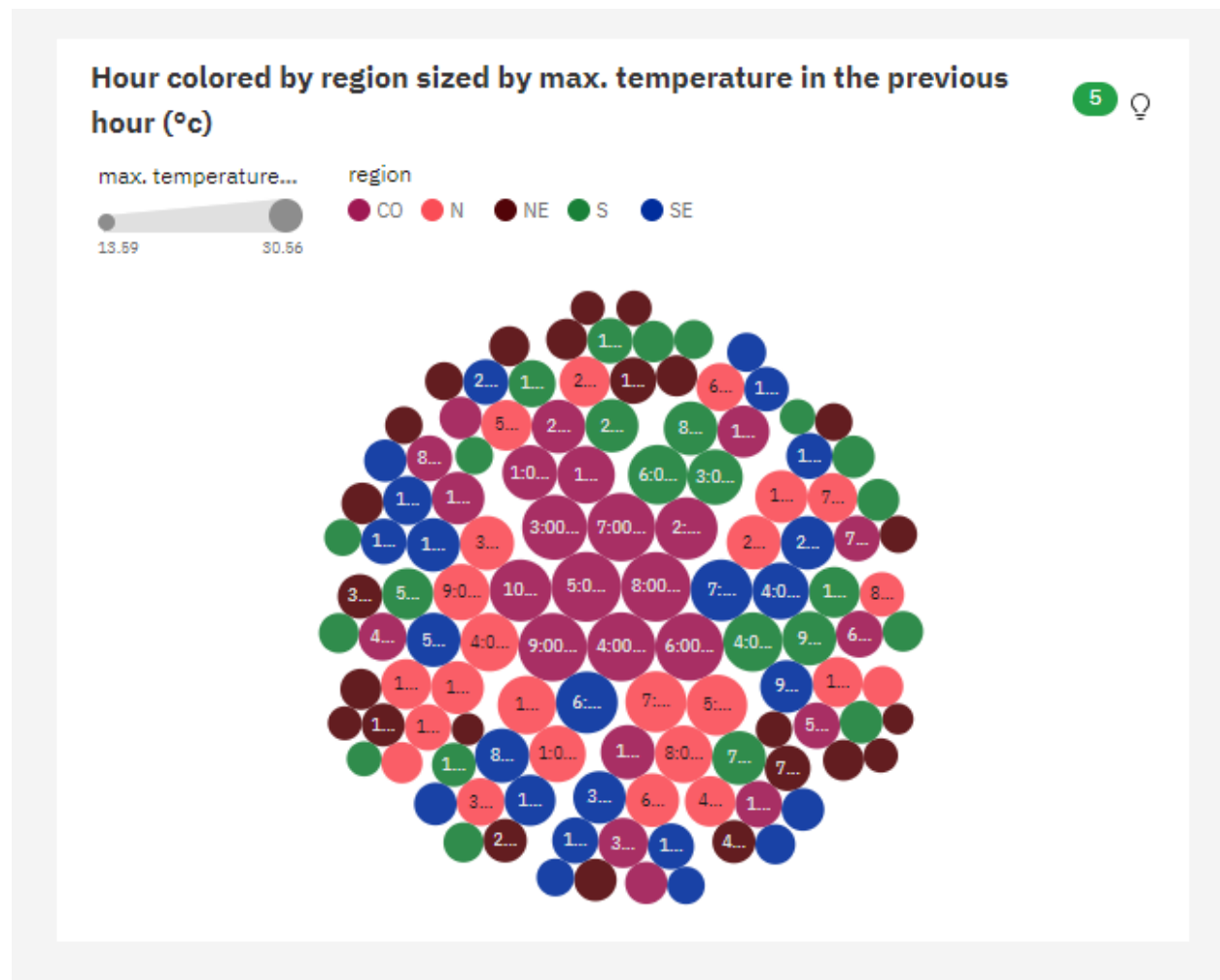
Visual 3 – Dew temperature in the previous hour and wind direction by region

As the heading says, this bar graph depicts the minimum dew temperature in all the regions in the previous hour and it also depicts the average wind direction in all the regions. The average of wind direction for all values of region is 121.6. The value of wind direction is unusually low when region is NE. The average of minimum dew temperature in the previous hour of all values in all regions is 13.79.



Visual 4 – Hour colored by region sized by maximum temperature

This bubble diagram shows the information of hours differentiated by region and sized by maximum temperature. As we can observe, the maximum temperature has been mostly in the CO region and the lowest of maximum temperature has mostly been in NE region.



Visual 5 – Air relative humidity by hour

The pie charts in the diagram below helps us to understand the relative humidity of each region differentiated by each hour from the dataset. Some of the facts we can get from these are that N and NE are the most frequently occurring regions with a combined total of 41.9% of humidity. On the other hand, the value of relative humidity is unusually low when it comes to CO region. Regarding the time, 13:00 is the most frequently occurring hour with a percentage of 4.7 of the total relative humidity.

