

VE472 HW4

吴佳遥 517370910257

Setup

```
1 CREATE TABLE dfs.tmp.`station` AS SELECT * FROM (  
2 select  
3 TRIM(SUBSTR(columns[0], 1, 11)) as id,  
4 TRIM(SUBSTR(columns[0],13,8)) as latitude,  
5 TRIM(SUBSTR(columns[0],22,9)) as longitude,  
6 TRIM(SUBSTR(columns[0],32,6)) as altitude,  
7 TRIM(SUBSTR(columns[0], 39, 2)) as state,  
8 TRIM(SUBSTR(columns[0], 42, 30)) as name  
9 from dfs.root.`/home/hadoop/ve472/weather/meta.csv`  
10 );
```

```
1 CREATE TABLE dfs.tmp.`weather` AS SELECT * FROM (  
2 select  
3 columns[0] as id,  
4 columns[1] as ob_date,  
5 columns[2] as ob_type,  
6 columns[3] as ob_value  
7 from dfs.root.`/home/hadoop/ve472/weather/2017.csv`  
8 );
```

```
1 CREATE TABLE dfs.tmp.`country` AS SELECT * FROM (  
2 select  
3 columns[0] as name,  
4 columns[1] as continent,  
5 columns[2] as fips  
6 from dfs.root.`/home/hadoop/ve472/weather/country_continent.csv`  
7 );
```

```
1 use dfs.tmp;
```

EX1.

1. Join operation

A `JOIN` clause is used to combine rows from two or more tables, based on a related column between them.

- `INNER JOIN`: selects records that have matching values in both tables.
- The `FULL OUTER JOIN` keyword returns all records when there is a match in left (table1) or right (table2) table records.

- The `LEFT JOIN` keyword returns all records from the left table (table1), and the matching records from the right table (table2). The result is 0 records from the right side, if there is no match.
- The `RIGHT JOIN` keyword returns all records from the right table (table2), and the matching records from the left table (table1). The result is 0 records from the left side, if there is no match.

2. Aggregate

An aggregate function is a function where the values of multiple rows are grouped together to form a single summary value.

The `COUNT()` function returns the number of rows that matches a specified criterion.

The `AVG()` function returns the average value of a numeric column.

The `SUM()` function returns the total sum of a numeric column.

The `MIN()` function returns the smallest value of the selected column.

The `MAX()` function returns the largest value of the selected column.

3. Advanced Nested Queries

1. The top five stations with the lowest average daily temperature

```
1 select station.name,weather.ob_value from weather
2 inner join station using (id)
3 where weather.ob_type = 'TAVG'
4 and LENGTH(weather.ob_value) > 0
5 and LENGTH(station.state) > 0
6 order by CAST(weather.ob_value as INTEGER) limit 5;
```

1	+-----+-----+	
2	name	ob_value
3	+-----+-----+	
4	Port Graham	-999
5	Monahan Flat	-999
6	Rocky Point	-999
7	RAM CREEK ALASKA	-733
8	RAM CREEK ALASKA	-733
9	+-----+-----+	

2. The top three stations with the highest max daily temperature in 20170831

```
1 select station.id,station.name,weather.ob_value from weather
2 inner join station using (id)
3 where weather.ob_type = 'TMAX'
4 and weather.ob_date = '20170831'
5 and length(weather.ob_value) > 0
6 and length(station.name) > 0
7 order by cast(weather.ob_value as float) desc limit 3;
```

```

1 | +-----+-----+-----+
2 | |   id   |  name  | ob_value |
3 | +-----+-----+-----+
4 | | KUM00040586 | JAHRA    | 484      |
5 | | IZ000040665 | KUT-AL-HAI | 484      |
6 | | IRM00040811 | AHWAZ     | 481      |
7 | +-----+-----+-----+

```

3. Min temperatures of stations with longitude between 29.5E and 30E

```

1 | select station.id,min(cast(weather.ob_value as float)) as tmin from weather
2 | inner join station using (id)
3 | where length(weather.ob_value) > 0
4 | and length(station.name) > 0
5 | and length(station.longitude) > 0
6 | and cast(station.longitude as float) < 30
7 | and cast(station.longitude as float) > 29.5
8 | group by station.id;

```

```

1 | +-----+-----+
2 | |   id   |  tmin  |
3 | +-----+-----+
4 | | ROM00015360 | -155.0 |
5 | | TUM00017155 | -143.0 |
6 | | SF001290070 | -107.0 |
7 | | FIE00146117 | 0.0    |
8 | | FIE00146598 | -260.0 |
9 | | FIE00144951 | -242.0 |
10 | | FIE00144887 | 0.0    |
11 | | FIE00144957 | -258.0 |
12 | | RSM00026268 | -277.0 |
13 | | FIE00144172 | 0.0    |
14 | | FIE00144877 | -288.0 |
15 | | NOE00133230 | 0.0    |
16 | | NOE00133210 | -180.0 |
17 | | FIE00144917 | 0.0    |
18 | | RSM00026167 | -260.0 |
19 | | BOM00033038 | -250.0 |
20 | | EGM00062318 | 0.0    |
21 | +-----+-----+

```

EX2

1. Perfect Weather

- Precipitation: $\leq 20mm$
- Average Temperature: $15 - 25^{\circ}C$
- Daily Temperature Amplitude: $\leq 7^{\circ}C$

2. Determine

First count days with perfect weather. See which stations have the most 10.

```

1 select tmax.id, COUNT(*) as perfect_days from weather tmax
2 inner join weather tmin on tmax.id = tmin.id and tmax.ob_date =
  tmin.ob_date and tmin.ob_type = 'TMIN'
3 inner join weather tavg on tmax.id = tavg.id and tmax.ob_date =
  tavg.ob_date and tavg.ob_type = 'TAVG'
4 inner join weather prcp on tmax.id = prcp.id and tmax.ob_date =
  prcp.ob_date and prcp.ob_type = 'PRCP'
5 where tmax.ob_type = 'TMAX'
6 and cast(prcp.ob_value as float) <= 150
7 and cast(tmax.ob_value as float) - cast(tmin.ob_value as float) <= 70
8 and cast(tavg.ob_value as integer) between 150 and 250
9 group by tmax.id
10 order by perfect_days desc
11 limit 10;

```

```

1 +-----+-----+
2 |      id      | perfect_days |
3 +-----+-----+
4 | SPE00120431 | 344          |
5 | SPE00120449 | 329          |
6 | SPE00120458 | 250          |
7 | SPE00120197 | 245          |
8 | ASN00009518 | 203          |
9 | ASN00009193 | 200          |
10 | SPE00120017 | 198          |
11 | USW00023188 | 187          |
12 | SP000060338 | 153          |
13 | MP000061995 | 140          |
14 +-----+-----+

```

Get the country

```

1 select * from country
2 where fips in (
3     select station.state from station
4     where station.id in (
5         select id from (
6             select tmax.id, COUNT(*) as perfect_days from weather tmax
7             inner join weather tmin on tmax.id = tmin.id and tmax.ob_date =
8             tmin.ob_date and tmin.ob_type = 'TMIN'
9             inner join weather tavg on tmax.id = tavg.id and tmax.ob_date =
10            tavg.ob_date and tavg.ob_type = 'TAVG'
11            inner join weather prcp on tmax.id = prcp.id and tmax.ob_date =
12            prcp.ob_date and prcp.ob_type = 'PRCP'
13            where tmax.ob_type = 'TMAX'
14            and cast(prcp.ob_value as float) <= 150
15            and cast(tmax.ob_value as float) - cast(tmin.ob_value as float)
16            <= 70
17            and cast(tavg.ob_value as integer) between 150 and 250
18            group by tmax.id
19            order by perfect_days desc
20        )
21        where perfect_days > 50
22    )
23 )
24 and length(continent) > 0;

```

```

1 | +-----+-----+-----+
2 | |   name   | continent | fips |
3 | +-----+-----+-----+
4 | | Suriname | SA       | NS   |
5 | | Turkmenistan | AS     | TX   |
6 | | Botswana | AF       | BC   |
7 | | Canada   | NA       | CA   |
8 | +-----+-----+-----+

```

My travel destination will be Suriname.

EX3

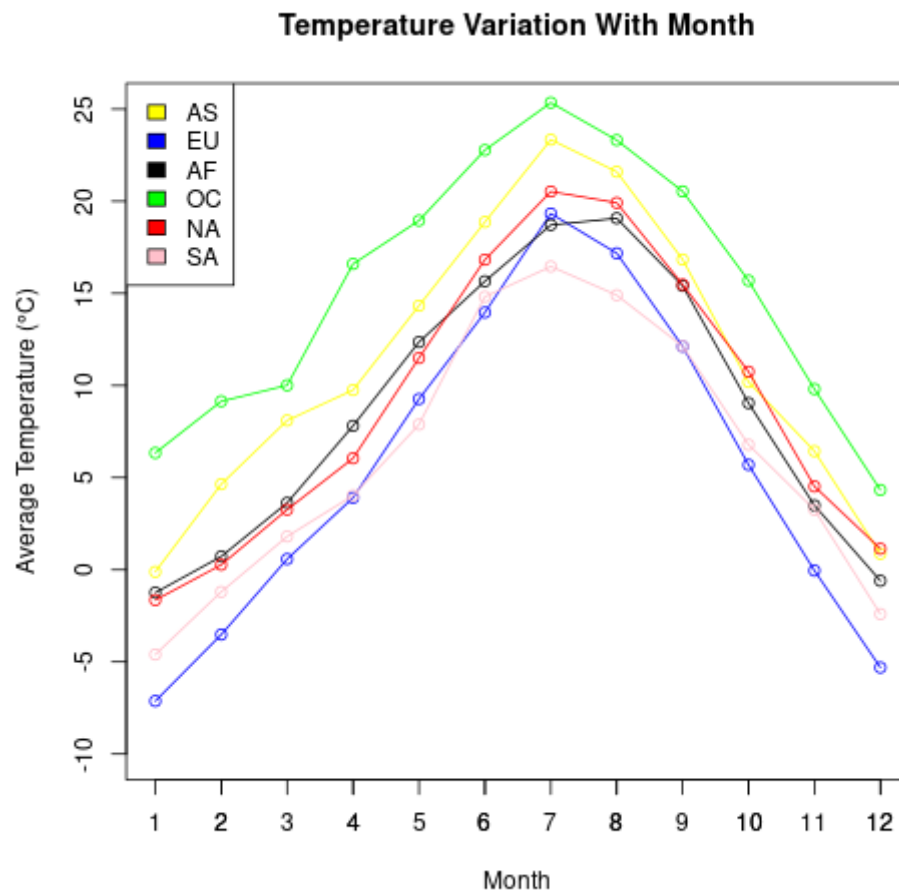
3.1

```

1 | select ob_month, avg(cast(ob_value as integer)) as avg_temperature
2 | from (
3 |     select *, substr(ob_date,5,2) as ob_month from weather
4 | )
5 | where id in (
6 | select st.id from station st
7 |     inner join country ctry on st.state = ctry.fips
8 |     where ctry.continent = 'AS'
9 | )
10 | and ob_type = 'TAVG'
11 | group by ob_month
12 | order by cast(ob_month as integer);

```

R script in `src/ex3.1.R`



3.2

```

1 select avg(cast(ob_value as integer)) as avg_temperature from weather
2 semicolon> where id in (
3     select st.id from station st
4     inner join country ctry on st.state = ctry.fips
5     where ctry.continent = 'SA'
6 )
7 and ob_type = 'TAVG';

```

R script in `src/ex3.2.R`

Annual average temperature (°C)

