# Exploratory Analysis of NOAA Weather Data

Benj, Grace and Mark

# Outline

# Dataset

# Dataset

Daily temperature (°C) and amount of precipitation (mm)

88 weather stations across the Philippines

January 1, 1960 to June 21, 2015

20261 rows and 89 columns

Longitude and latitude of the stations

# Dataset

National Oceanic and Atmospheric Administration's (NOAA) Integrated Surface Data (ISD)

ftp://ftp.ncdc.noaa.gov/pub/data/noaa/

1901 to 2016

293 Countries

# Data Cleaning

# Data Cleaning

88 weather stations across the Philippines but reduced to 85 stations

- Baler + Baler Radar
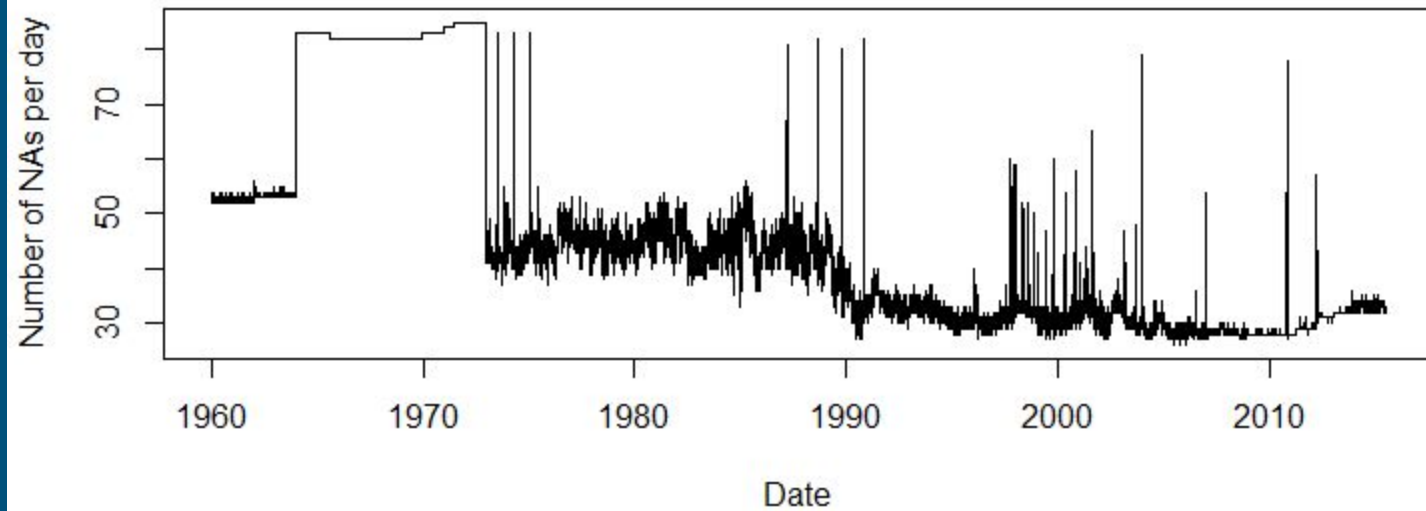- Basco + Basco Radar
- Davao + Davao Airport

# Missing Values

There's a lot of missing (NA) values in the dataset probably caused by

- Station is not yet established
- Station experienced difficulties in reading data
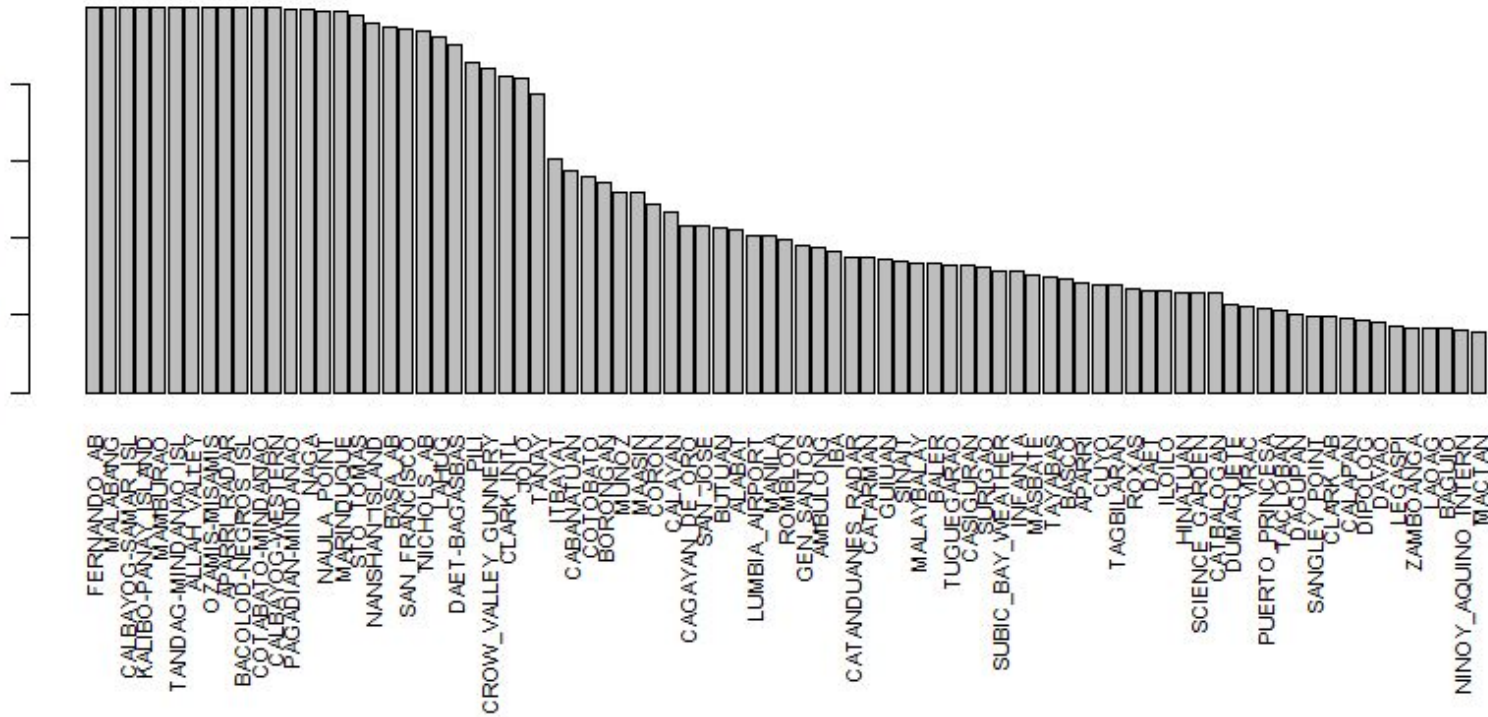- Etc

The following plots are:

- # of NAs over a period of time
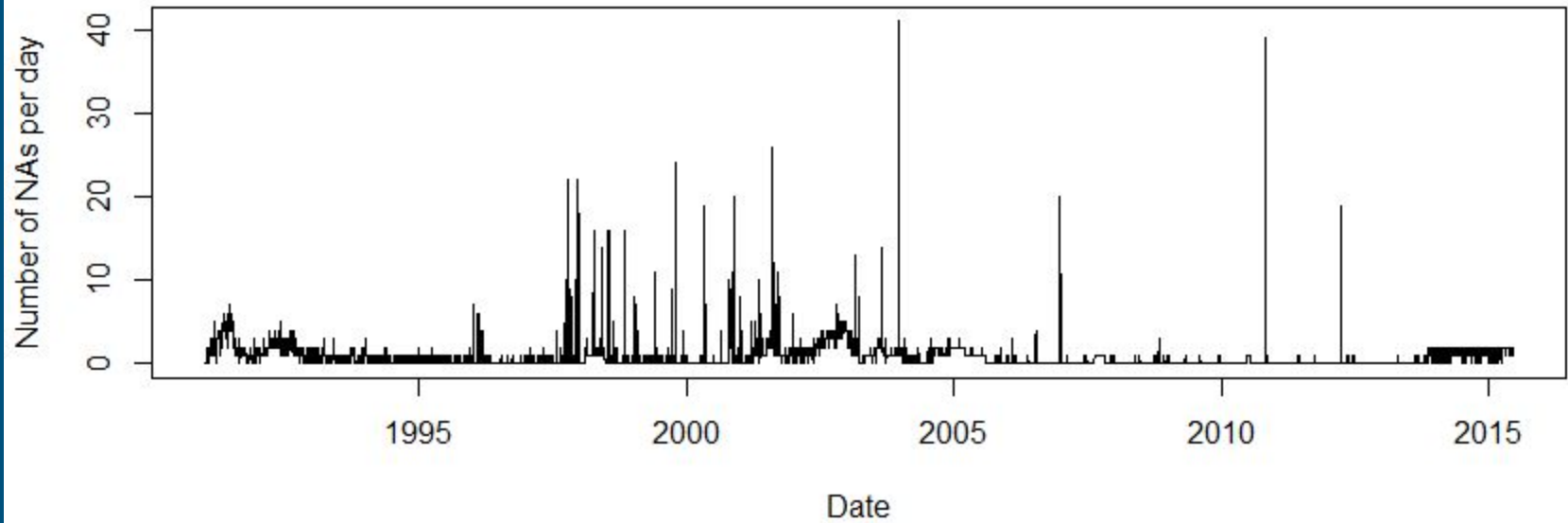- Distribution of NAs among the stations
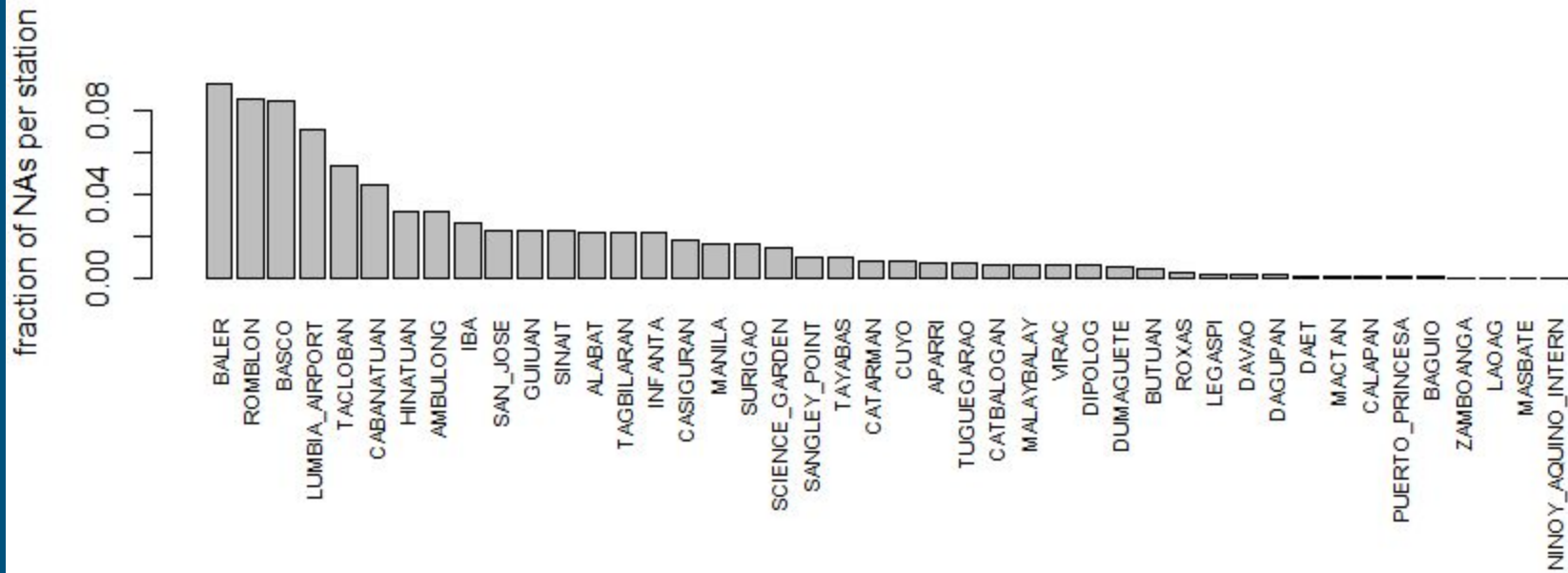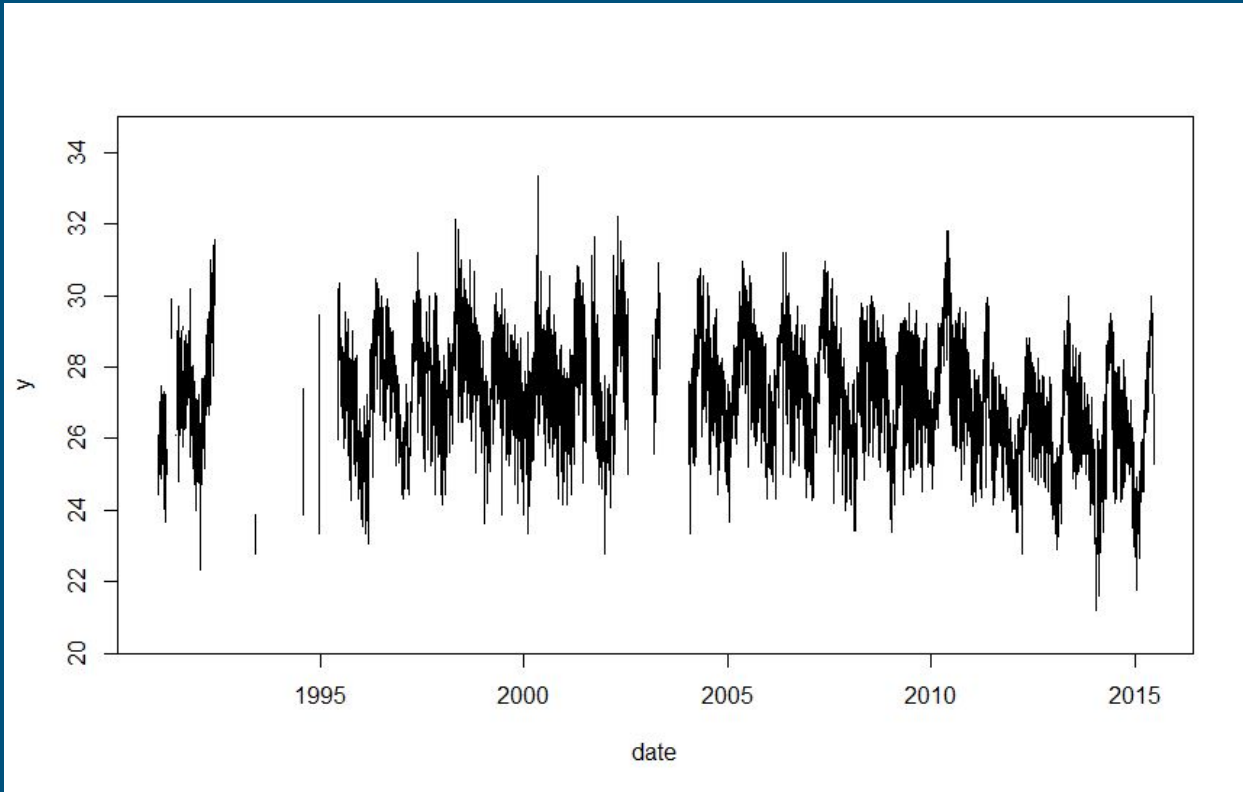
A lot of missing data between 1966 to 1974

Stations with the most number of NAs

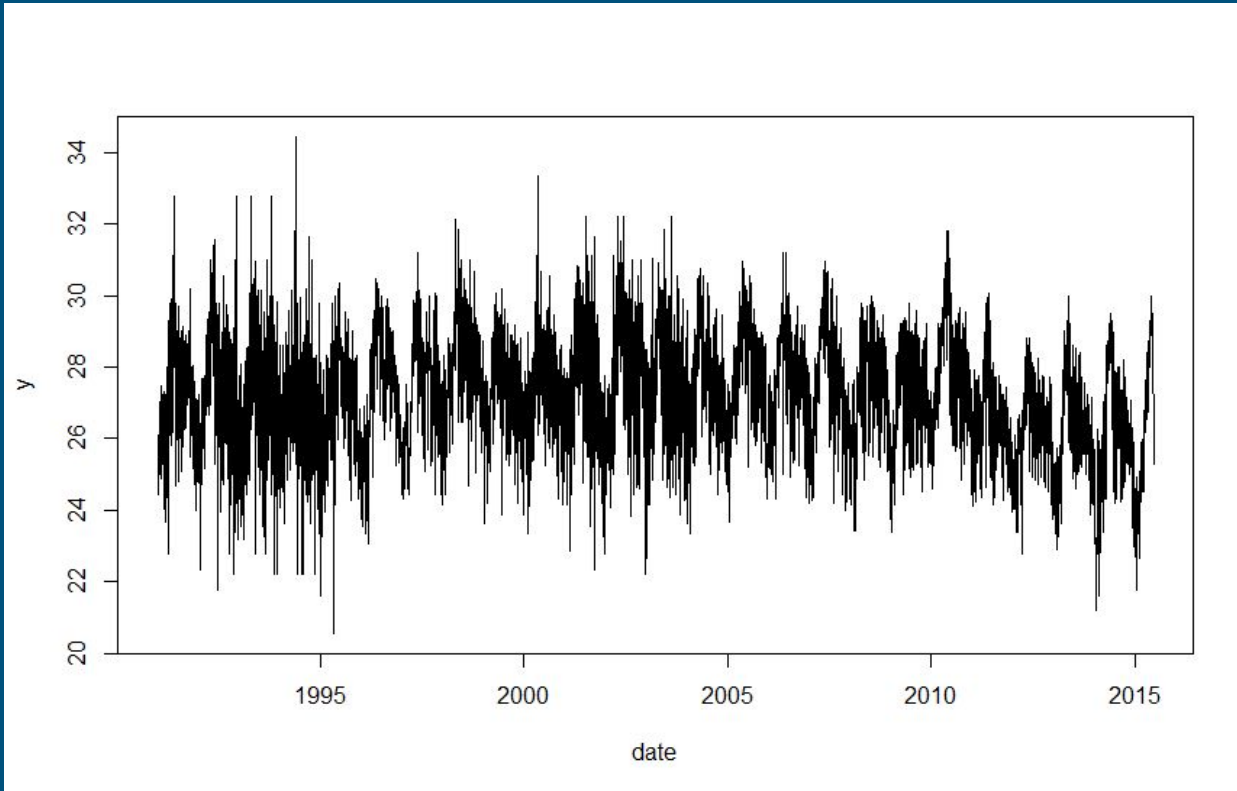We limit the timeframe between 1990 to 2015

Removed some of the stations

Data with missing values

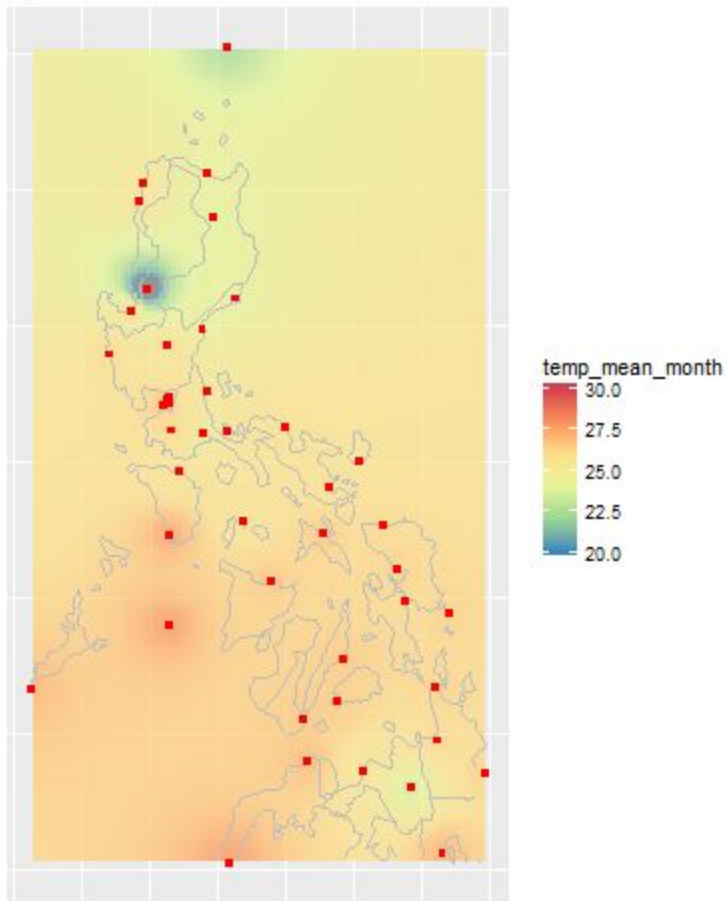Data with missing values with replaced with mean

Data with missing values with replaced with **predictive mean matching**
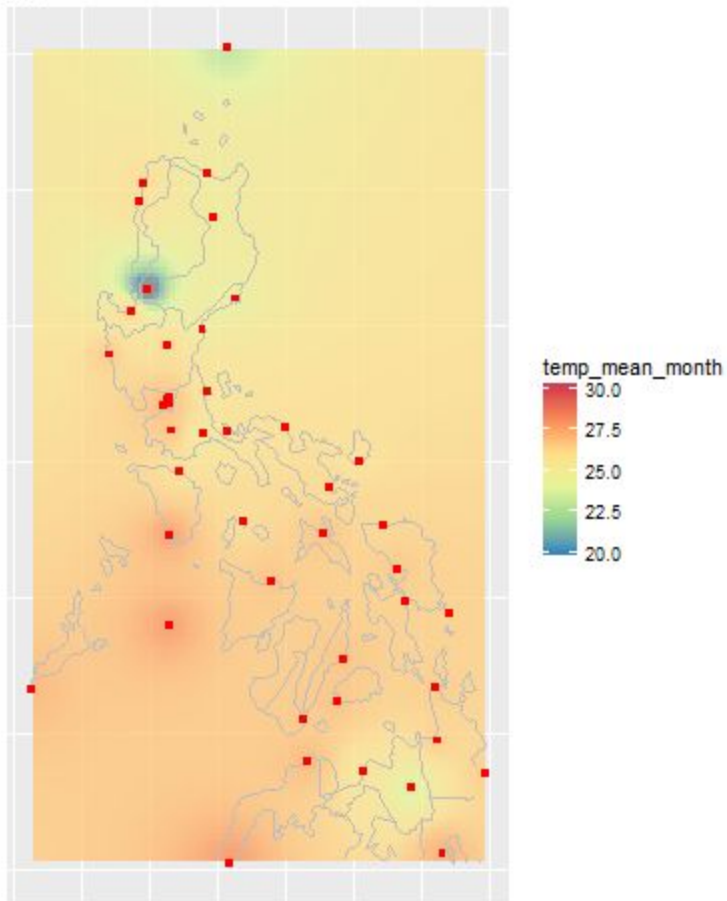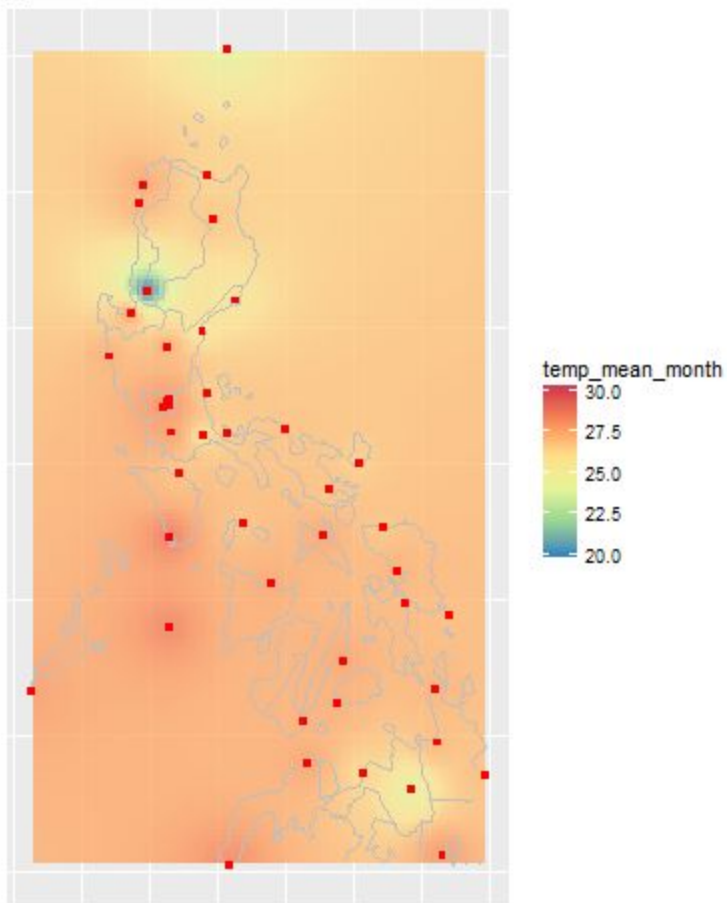
# Visualizations
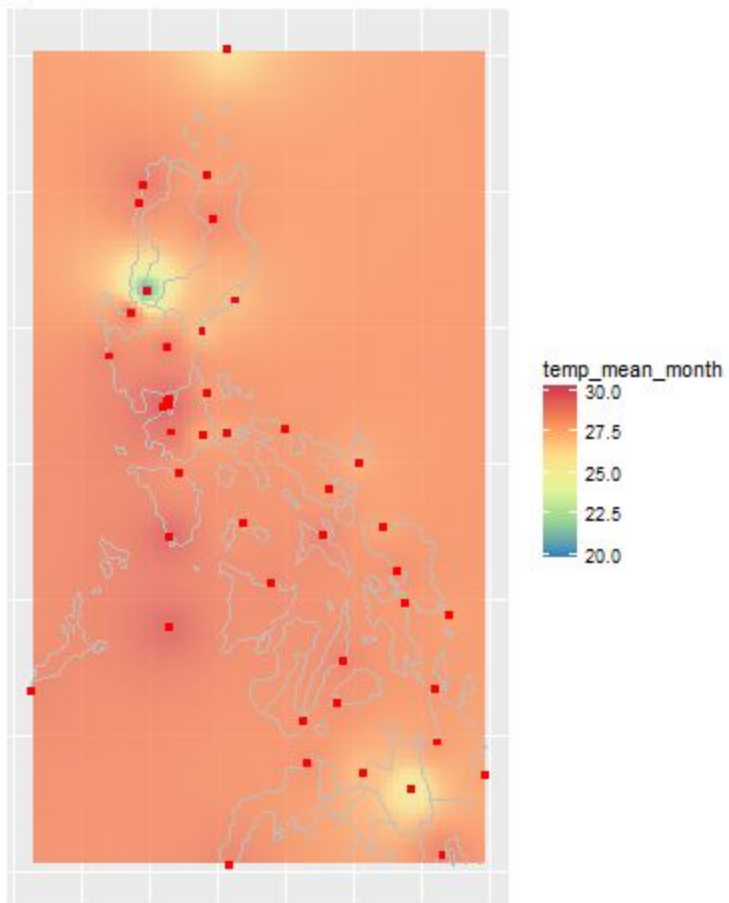
# Temperature

Average Monthly Temperature (°C)

January

Average Monthly Temperature (°C)

February

03

Average Monthly Temperature (°C)

March

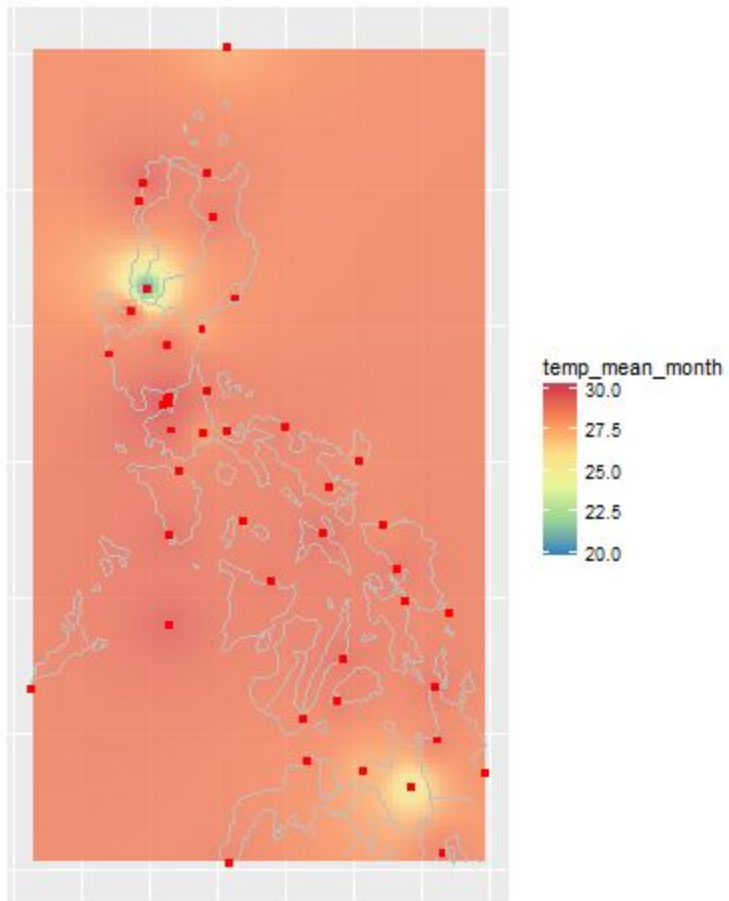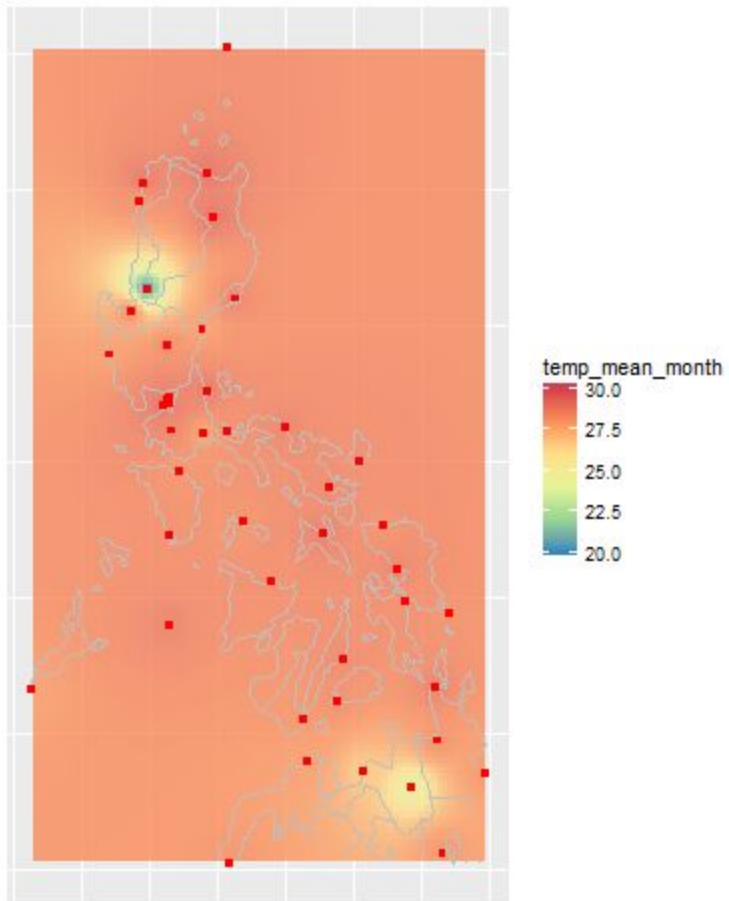**Average Monthly Temperature (°C)**
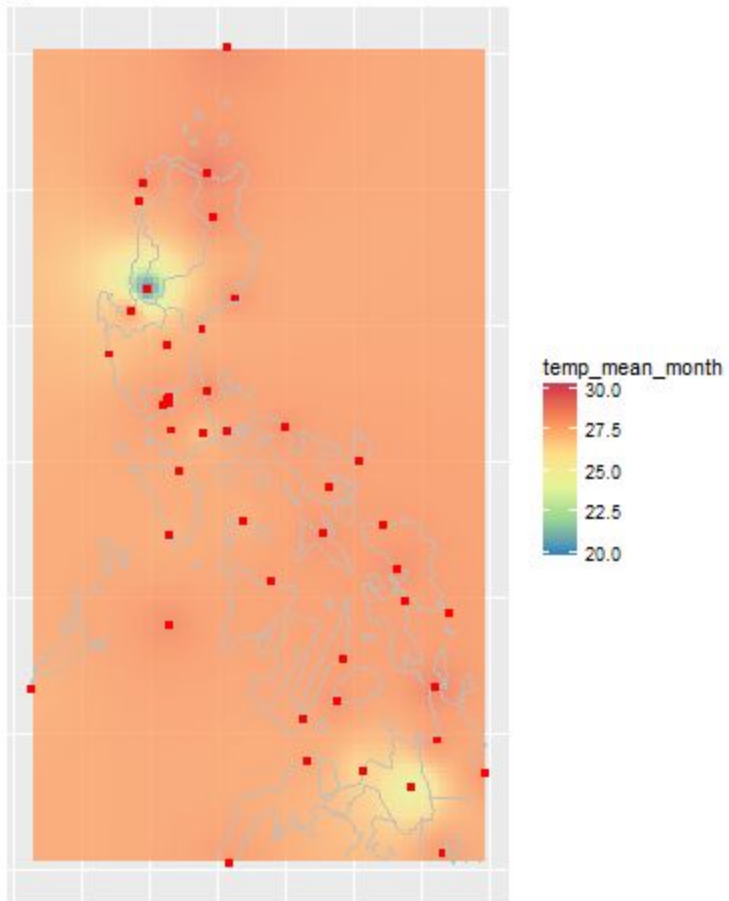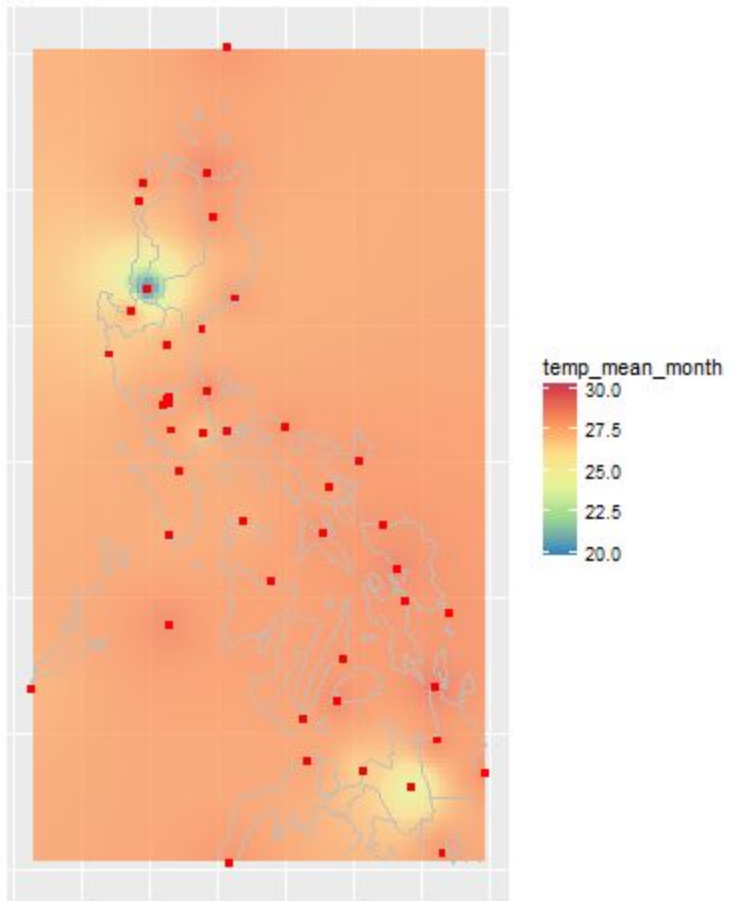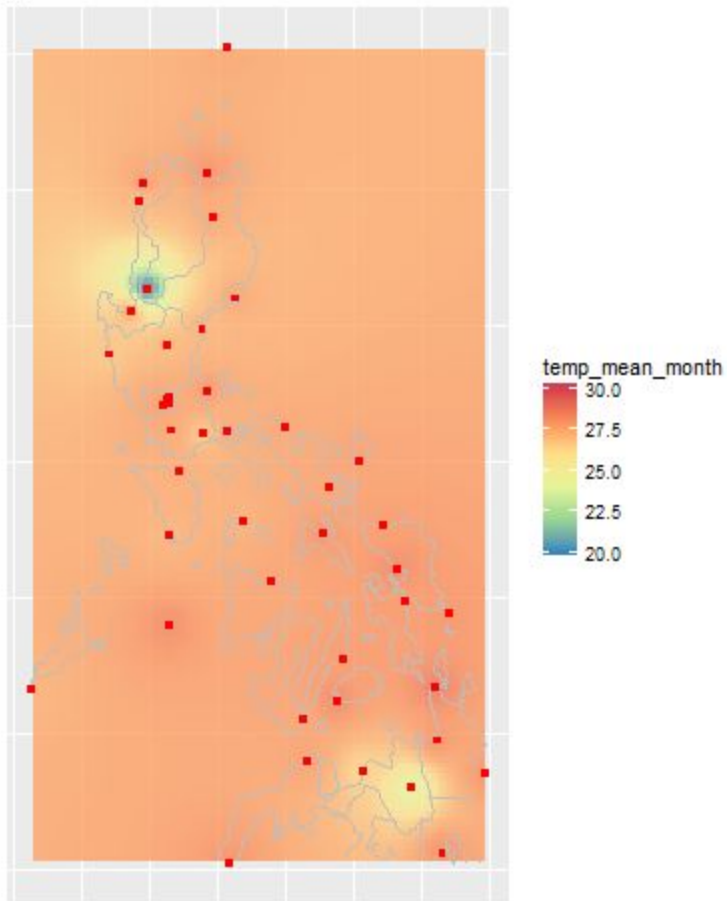
**May**

Average Monthly Temperature (°C)

July

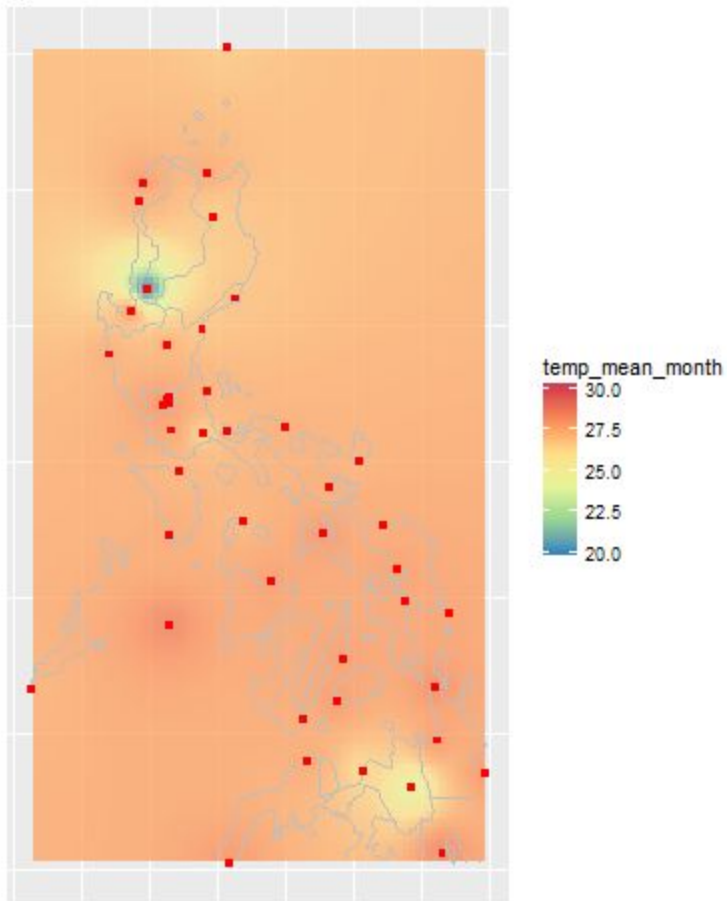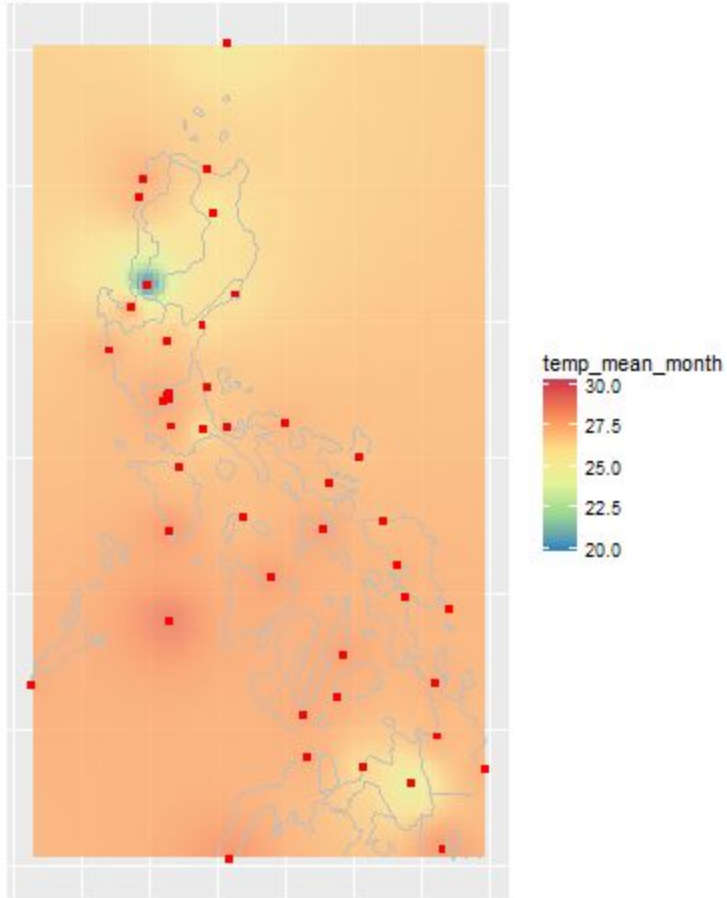Average Monthly Temperature (°C)

August

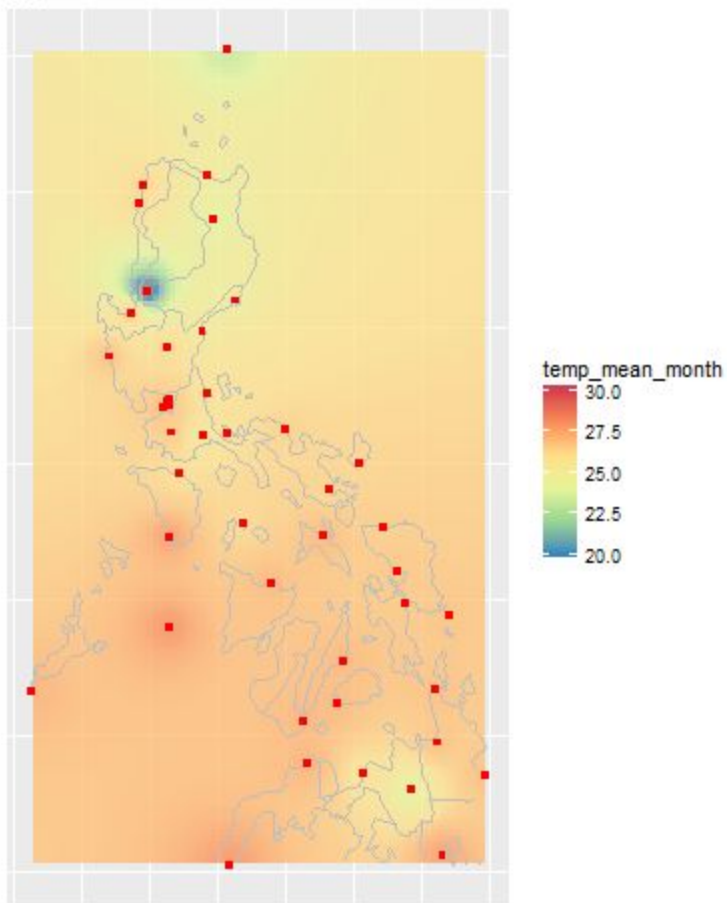Average Monthly Temperature (°C)

September

Average Monthly Temperature (°C)
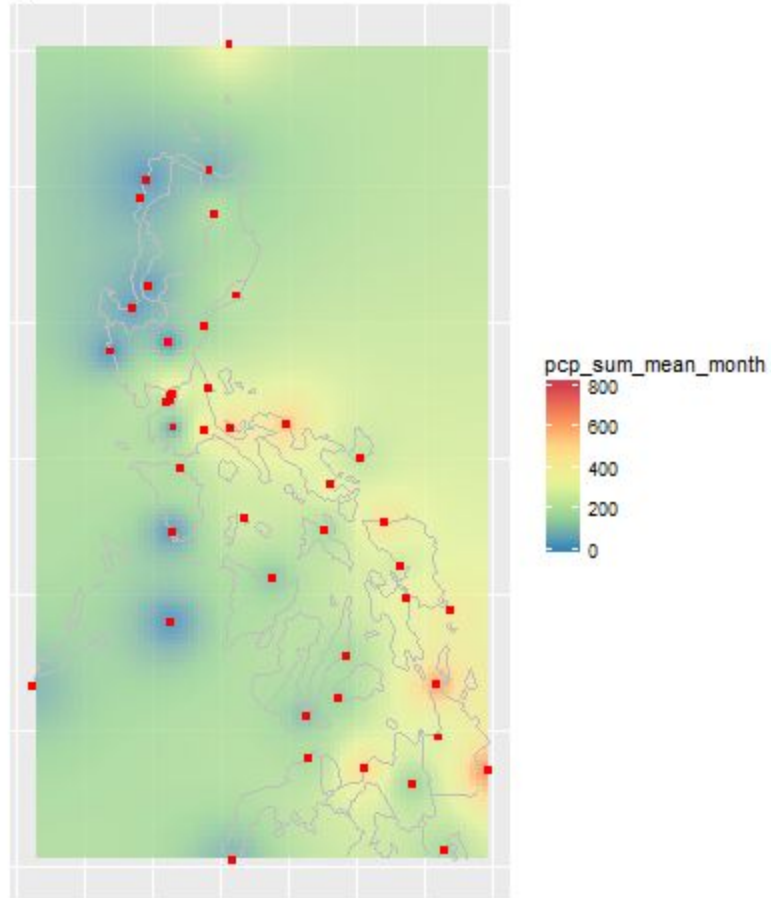
October

**Average Monthly Temperature (°C)**

# November

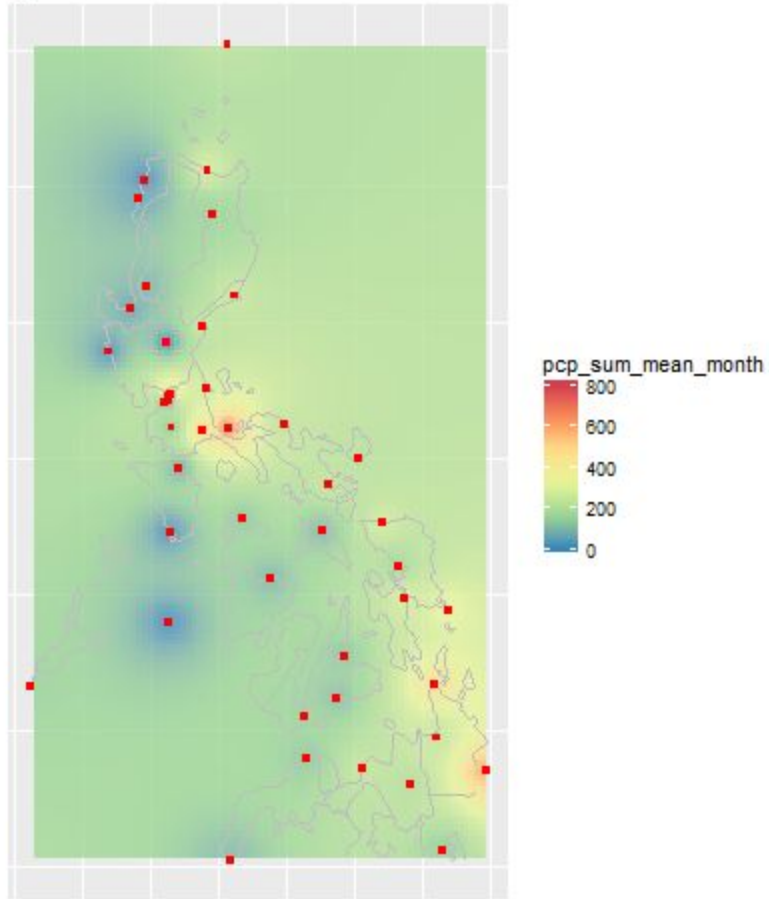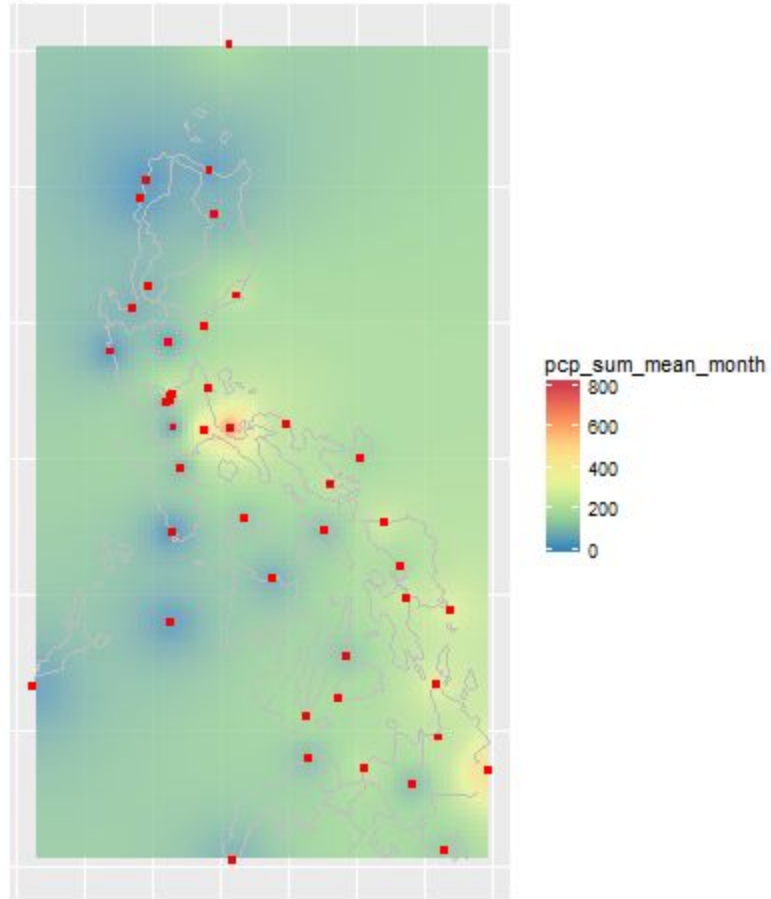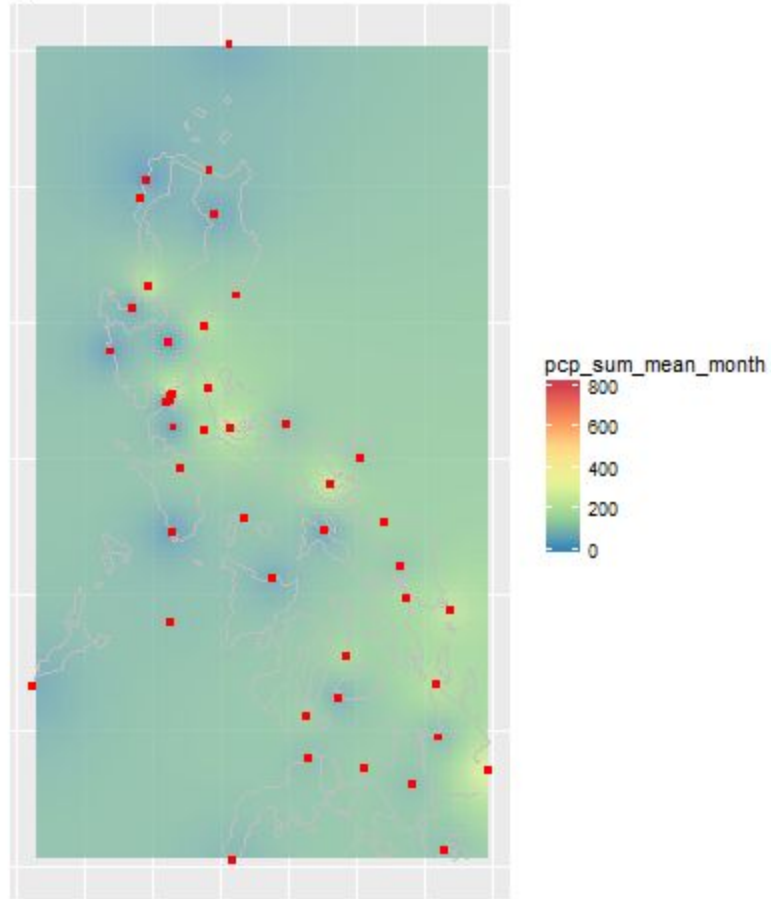Average Monthly Temperature (°C)

December

# Precipitation

**Average Monthly Precipitation (mm)**
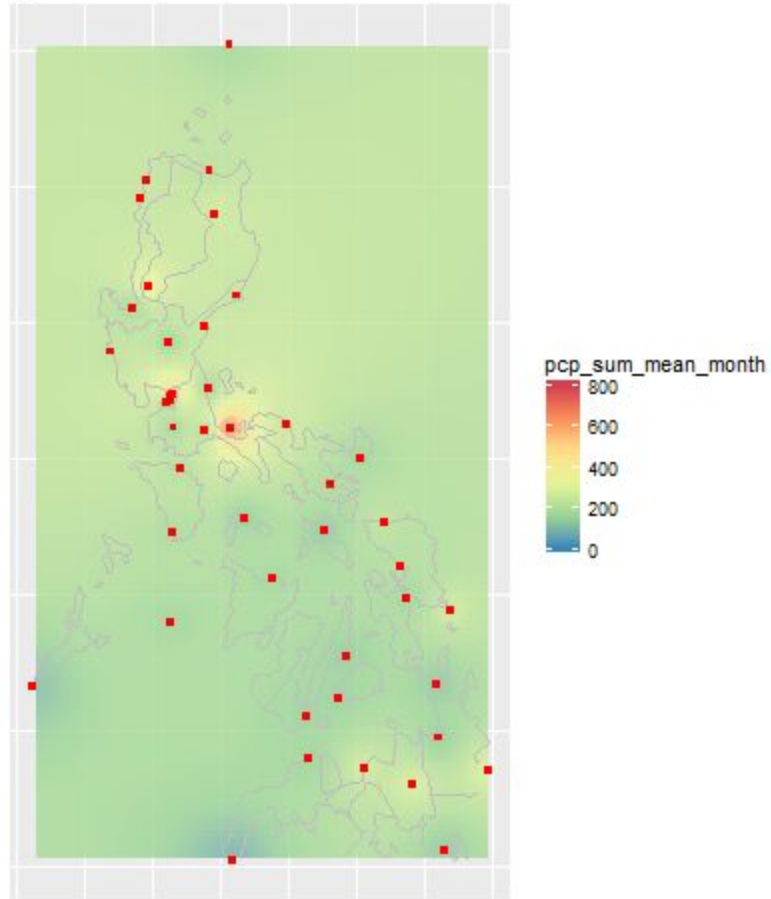
# January

Average Monthly Precipitation (mm)

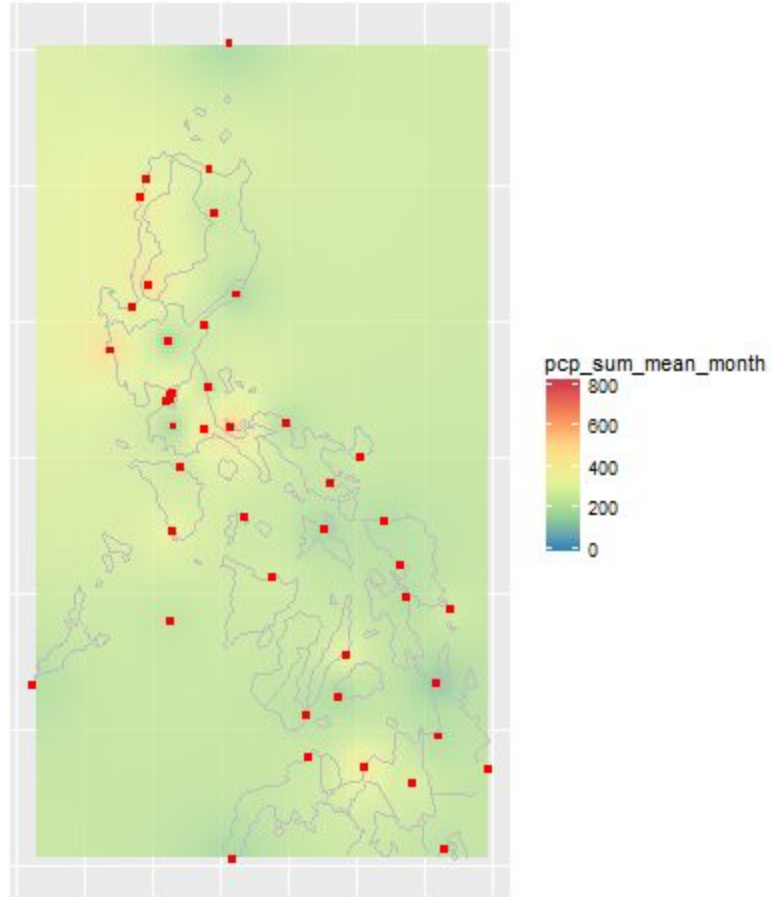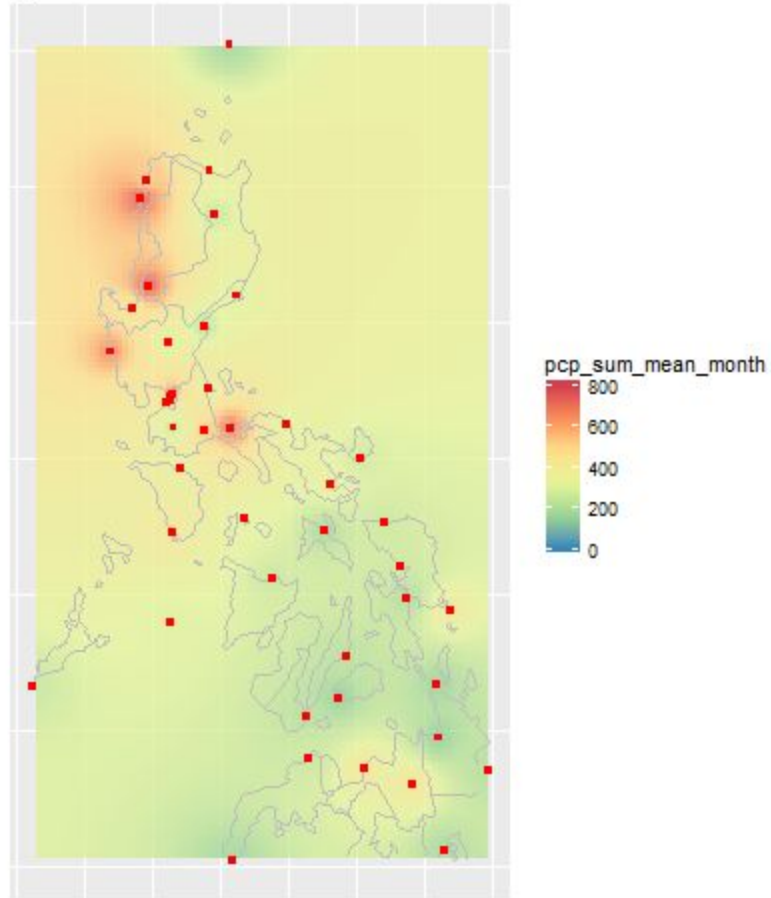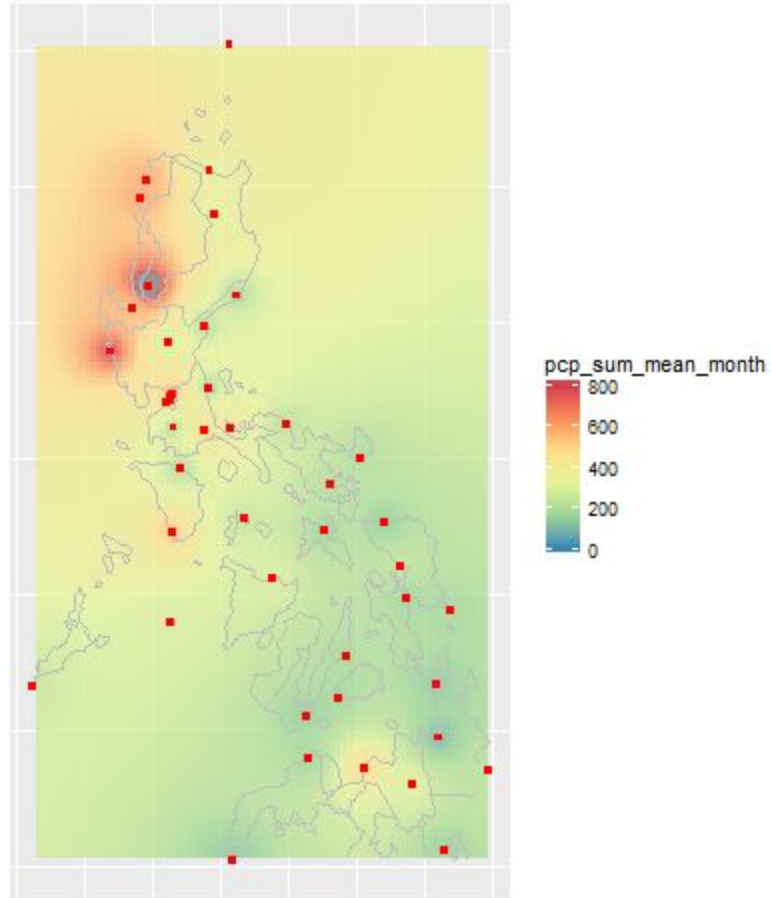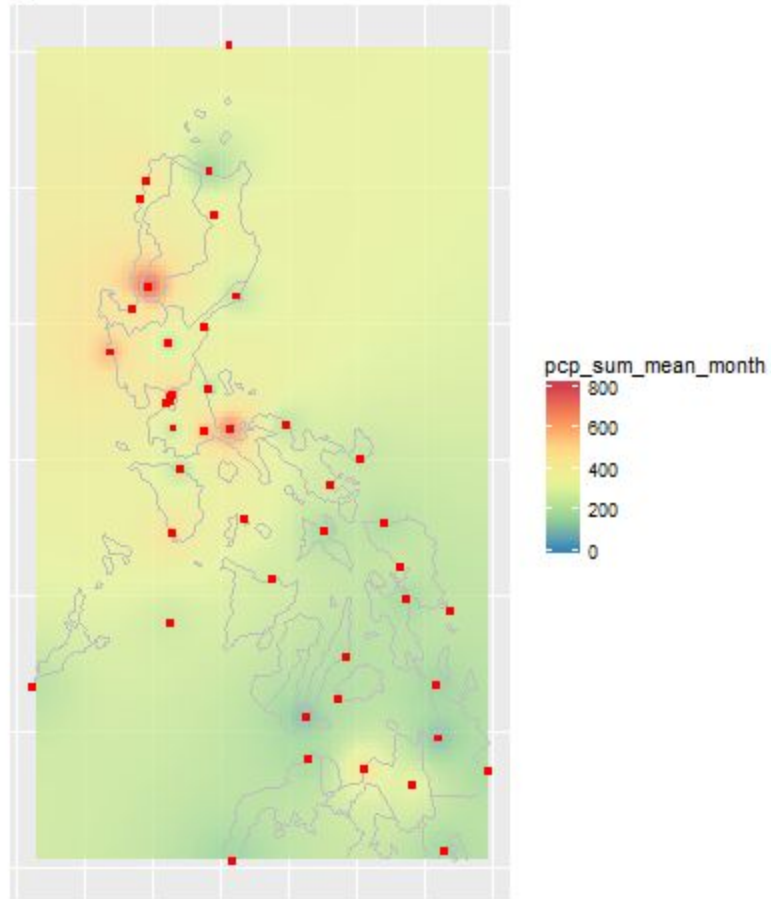March

**Average Monthly Precipitation (mm)**
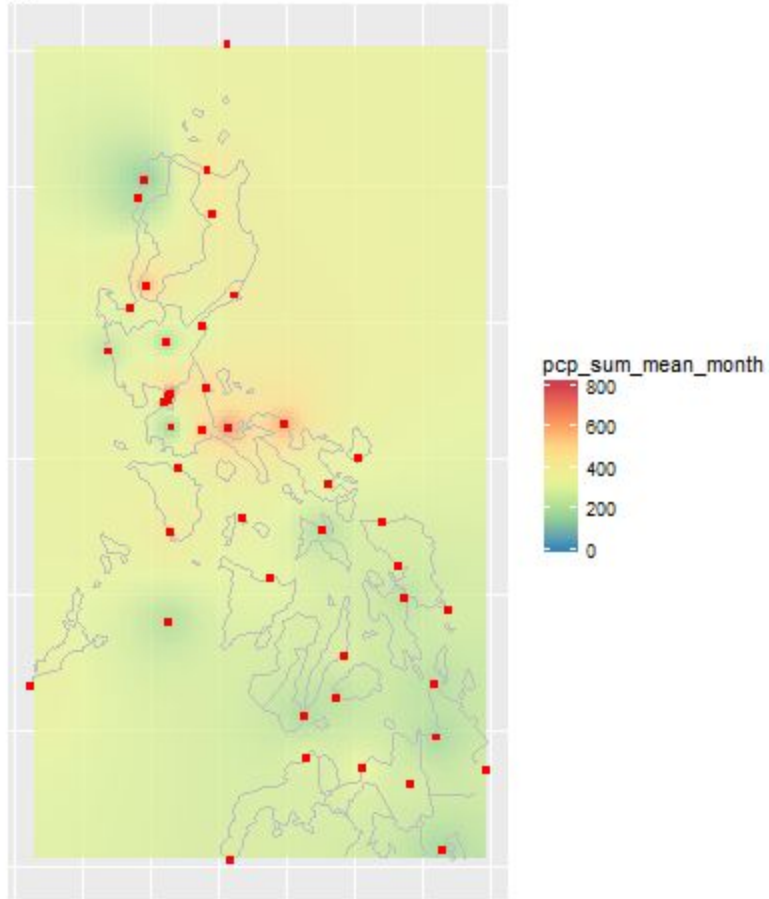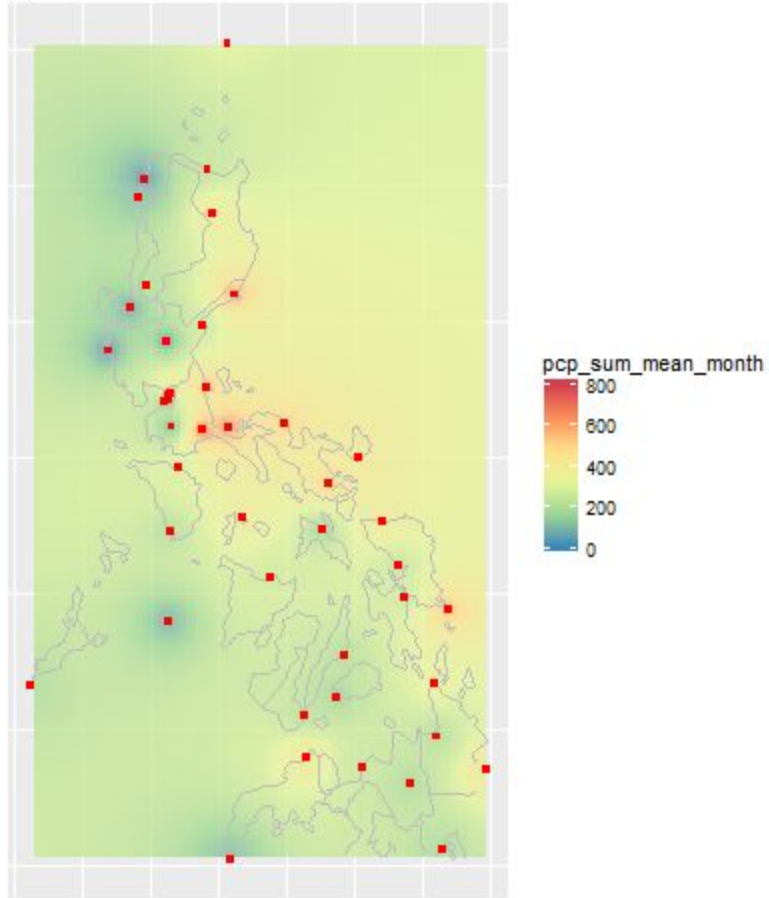
**May**

Average Monthly Precipitation (mm)

June
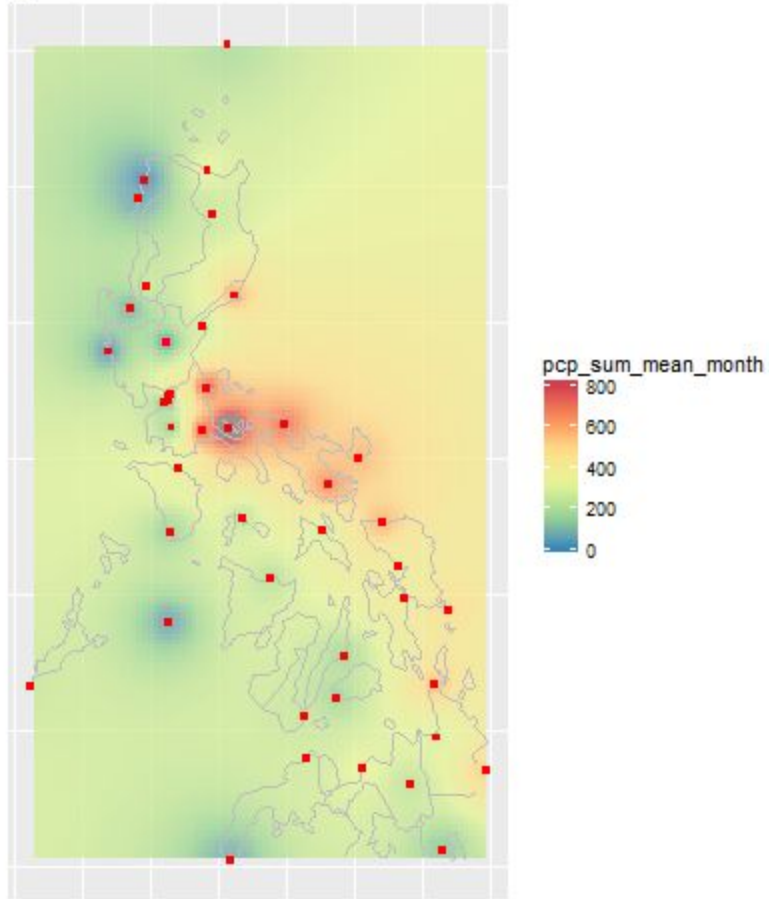
Average Monthly Precipitation (mm)

August

Average Monthly Precipitation (mm)
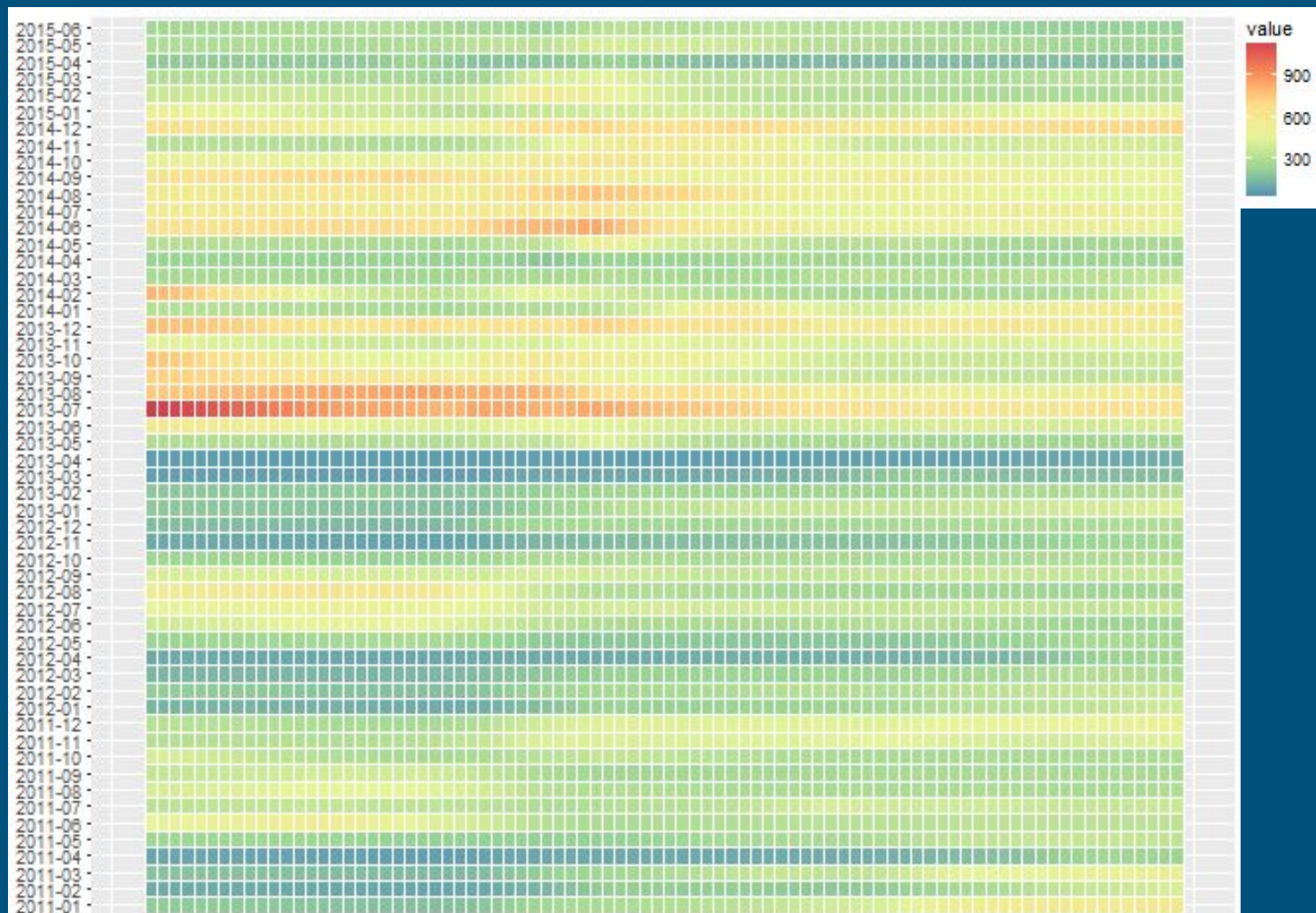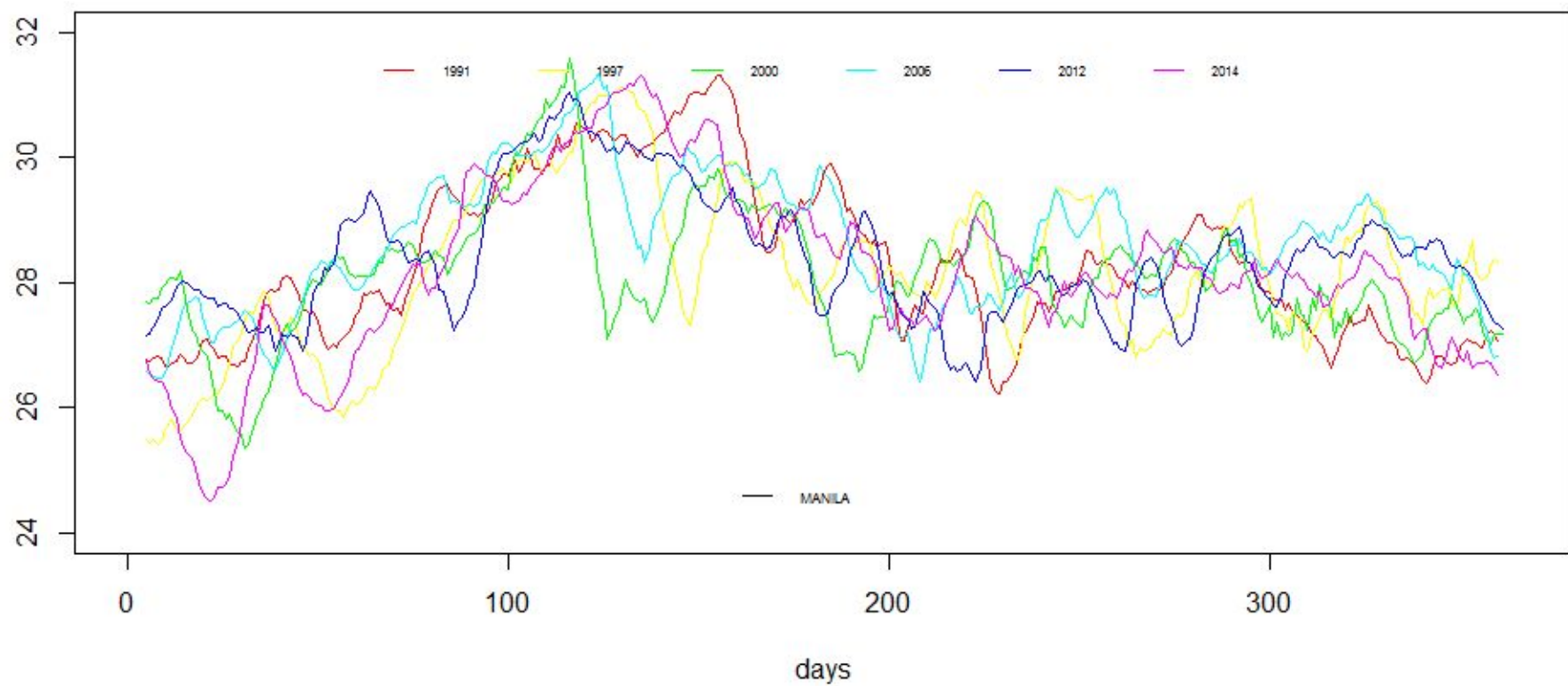
October

Average Monthly Precipitation (mm)
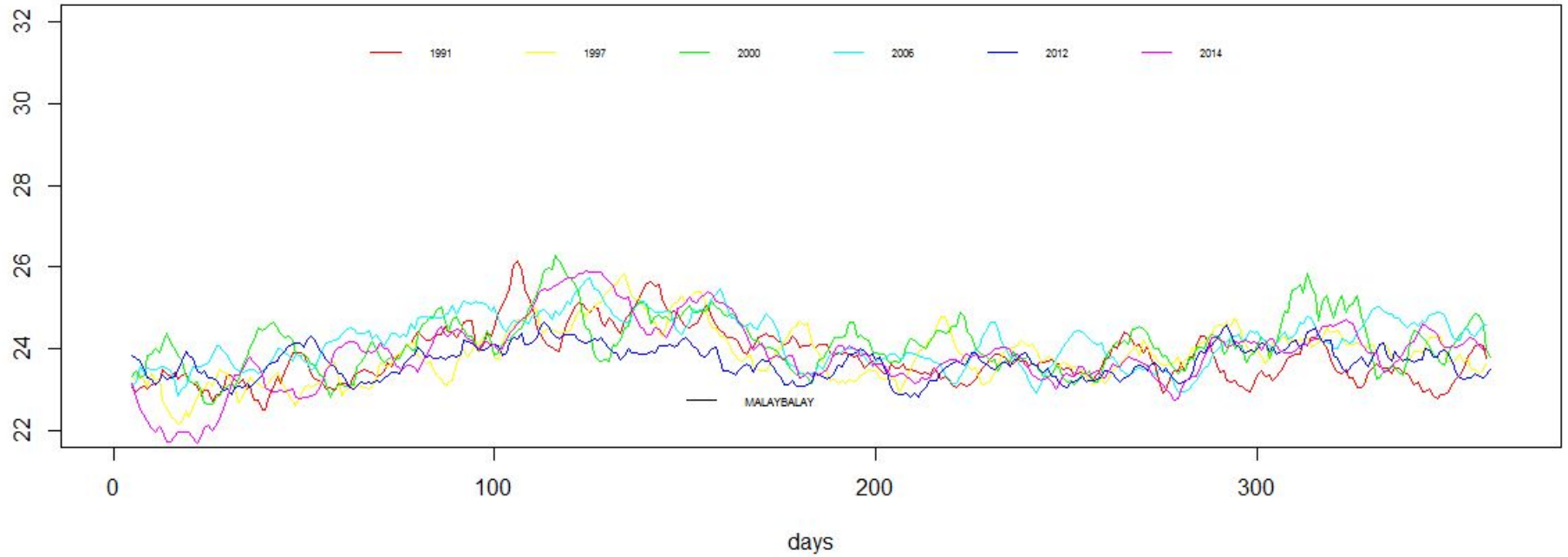
November

Average Monthly Precipitation (mm)
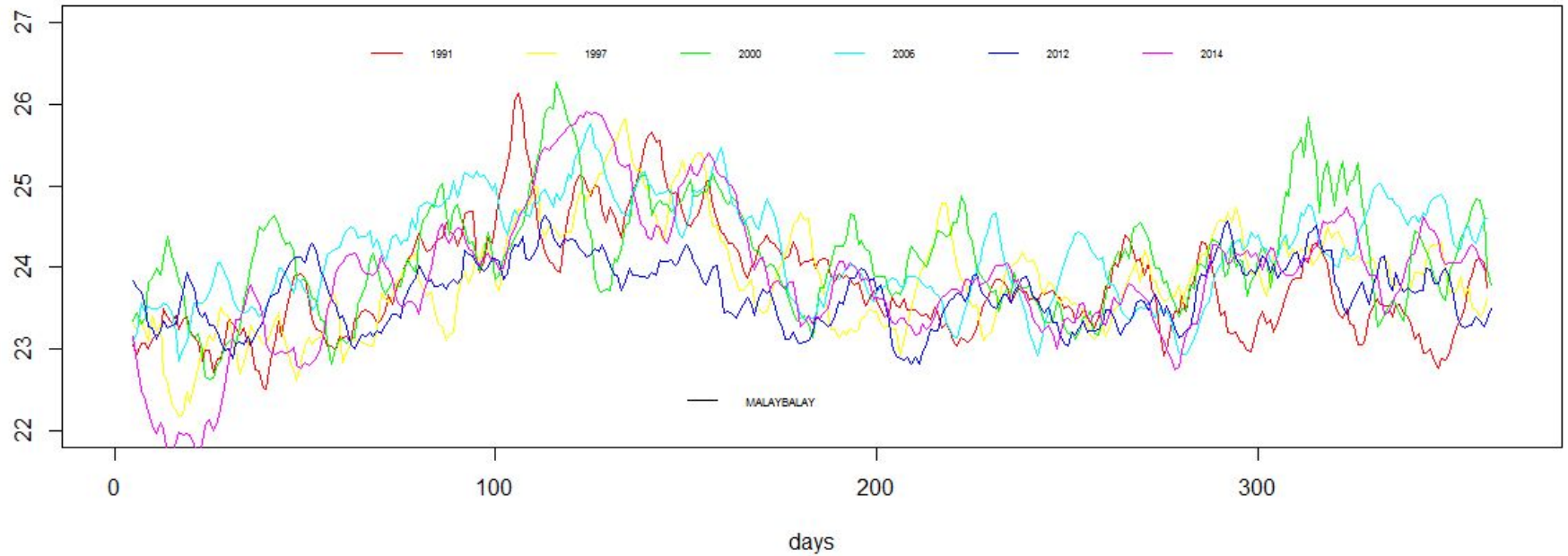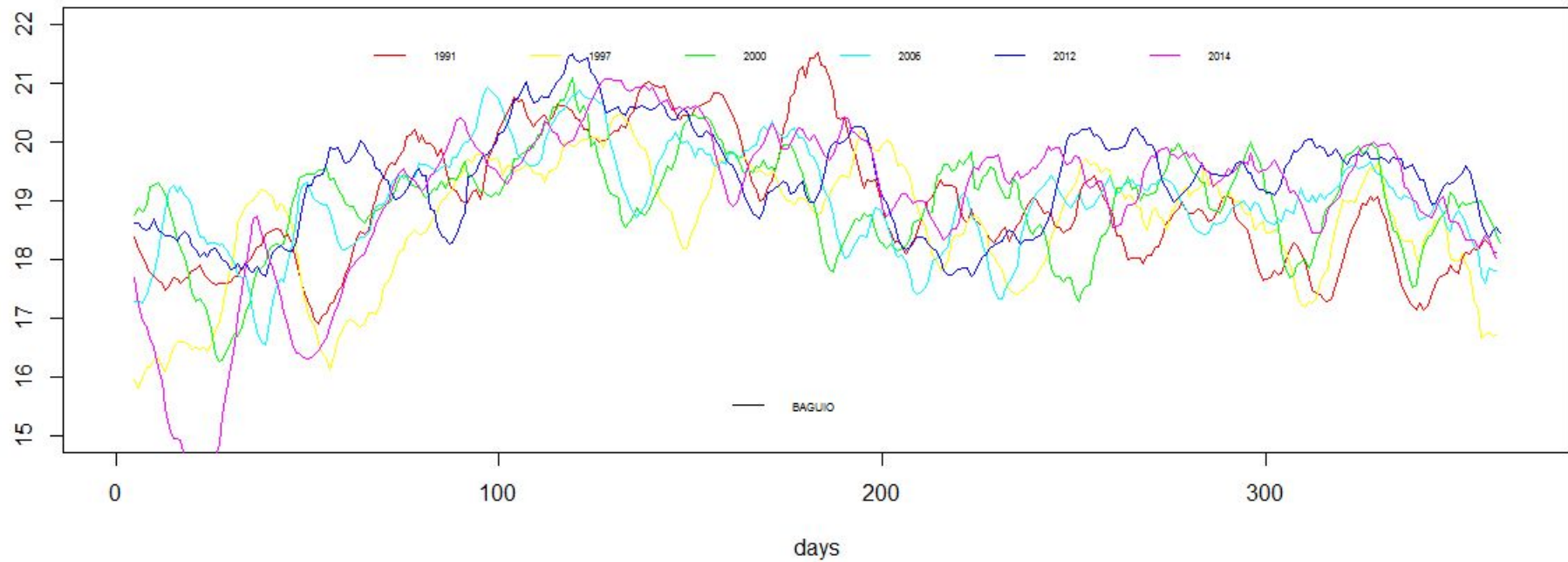
December

# Temperature Time Series

MANILA Temperature

MALAYBALAY Temperature

MALAYBALAY Temperature

# BAGUIO Temperature

# Change of temperature over the years

Legend: 1991, 1994, 1997, 2000, 2003, 2006, 2009, 2012, 2014

days

all

Temperature Correlation

Precipitation Correlation

# Variogram

We want to find out if there's:

1. Spatial continuity
2. Lag

**Variogram**

- variance vs distance
- variogram function from gstat
- 6000 samples

Variogram + fit.variogram

**Fit.variogram**

Fits the existing variogram to a model (e.g. gaussian, exponential)

# Temperature (semi)variogram

# Temperature (semi)variogram

# Temperature (semi)variogram

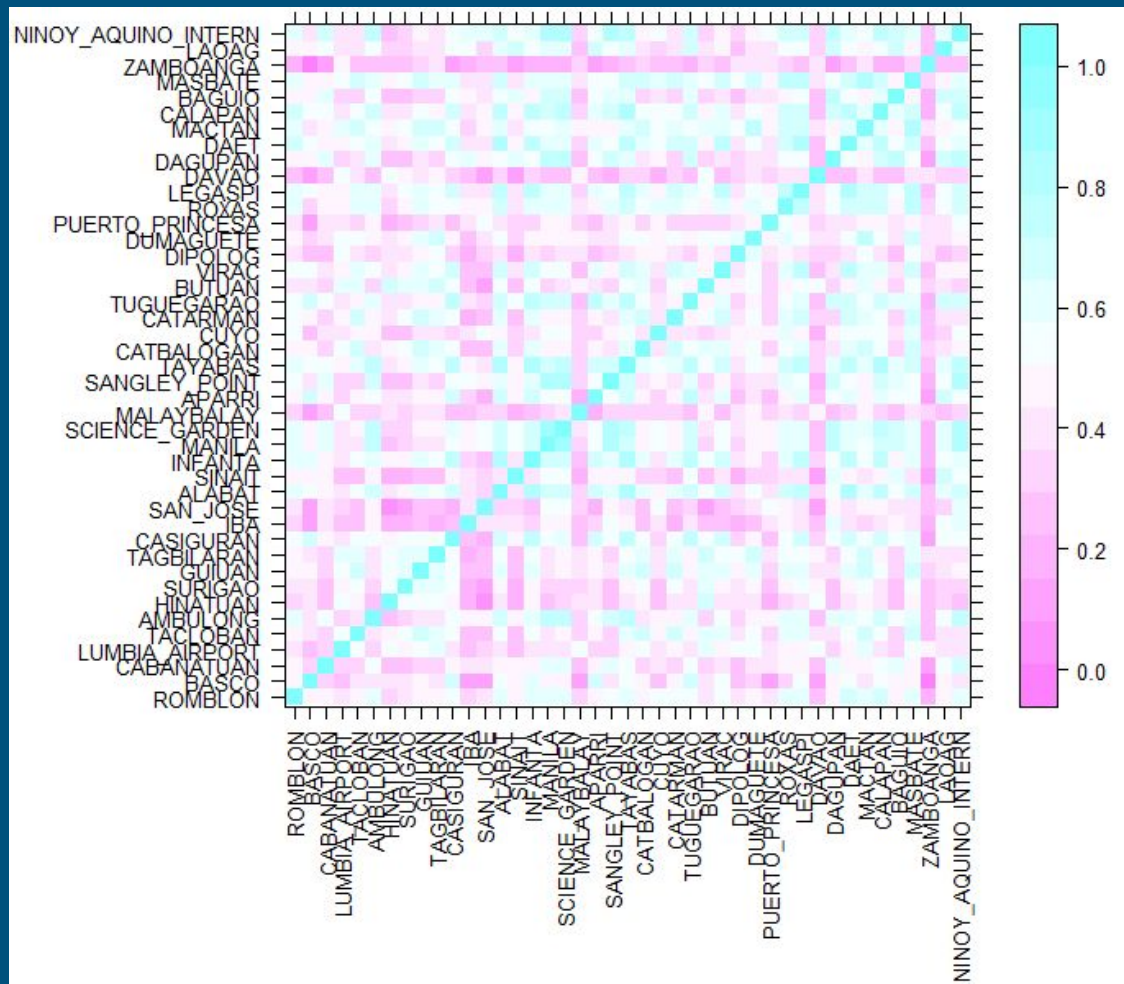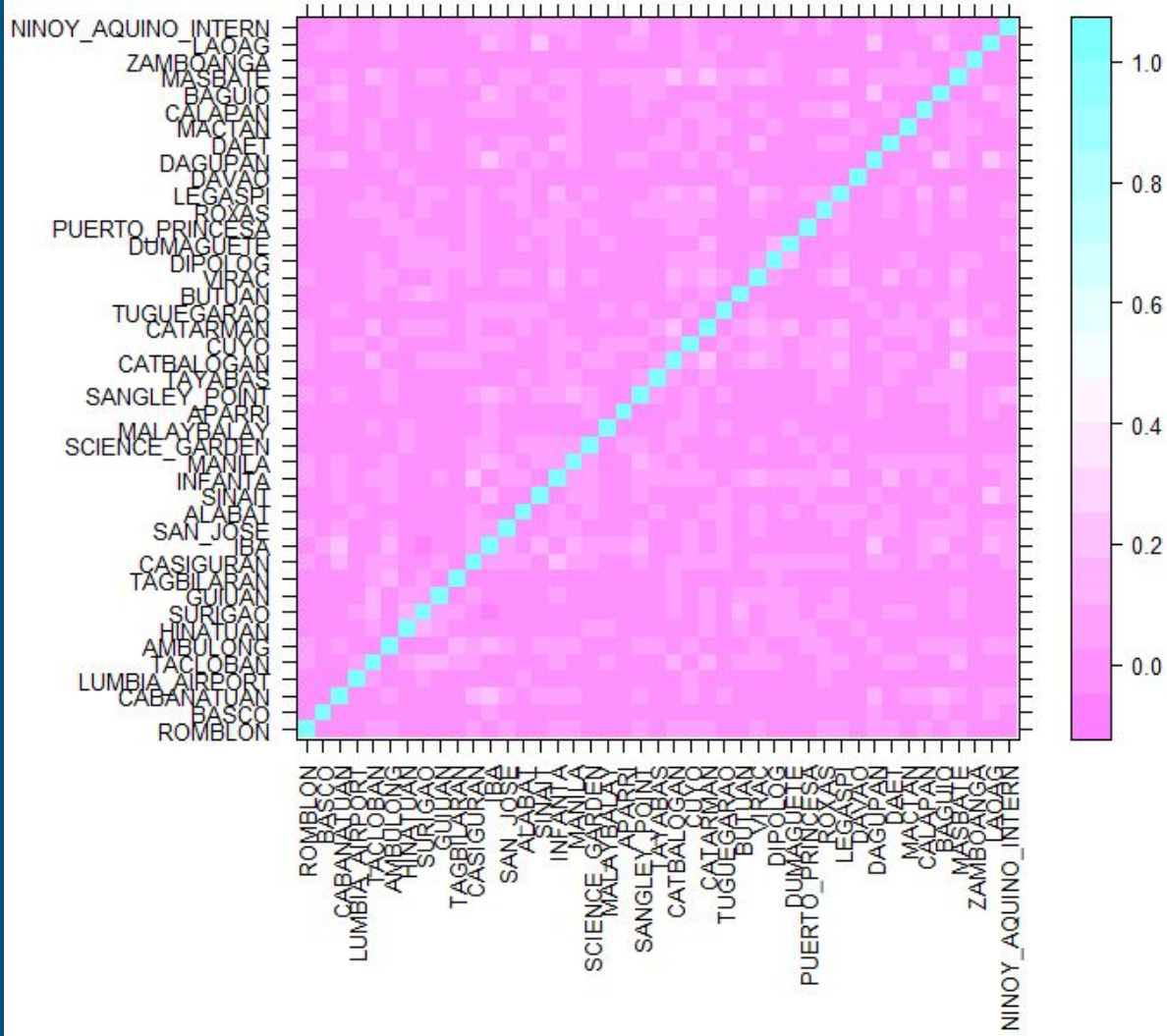Convergence error: gaussian, spherical and exponential models

Periodic model

Initial guess: no pattern or trend

➔ Range - indicates at which distance the variogram reaches the sill value (or where it levels off)
➔ Range = 0.00

# Precipitation

# Temperature Full Dataset

No convergence on
available models

# Cluster Analysis

# Hierarchical Clustering

"Closeness" of the stations
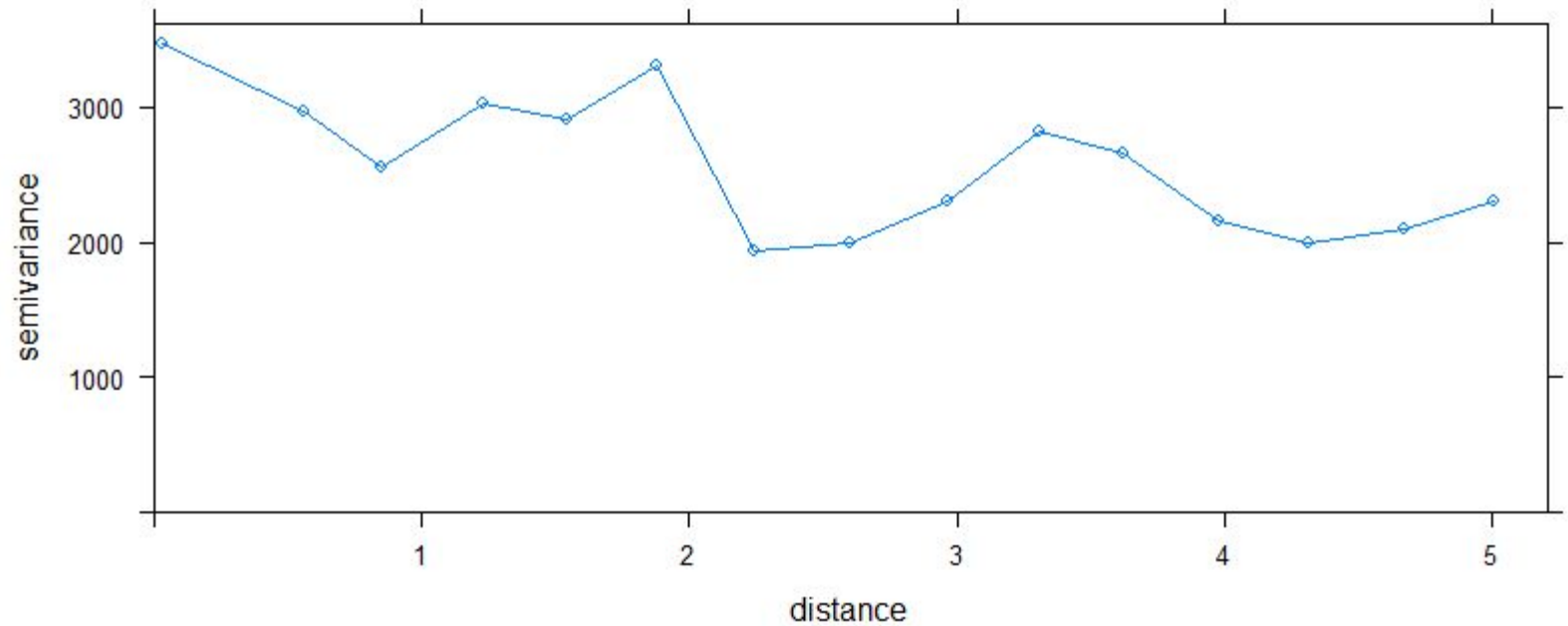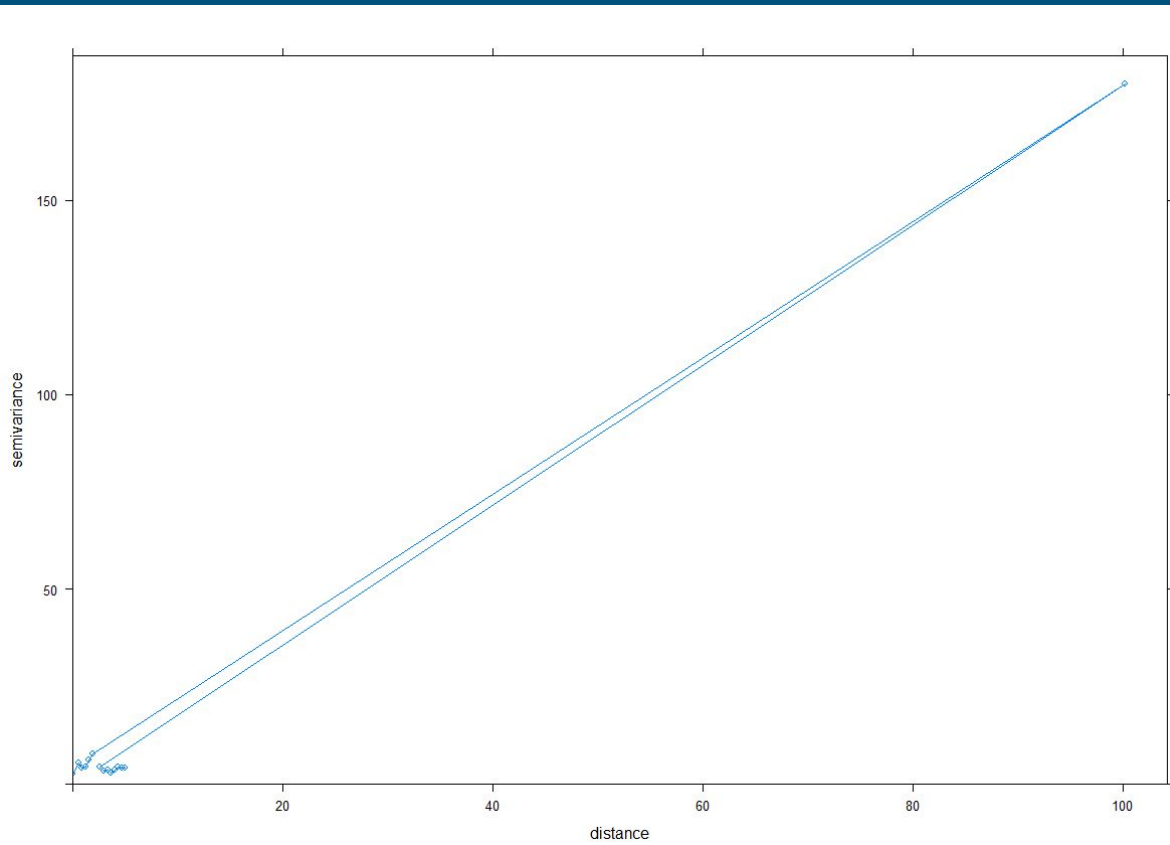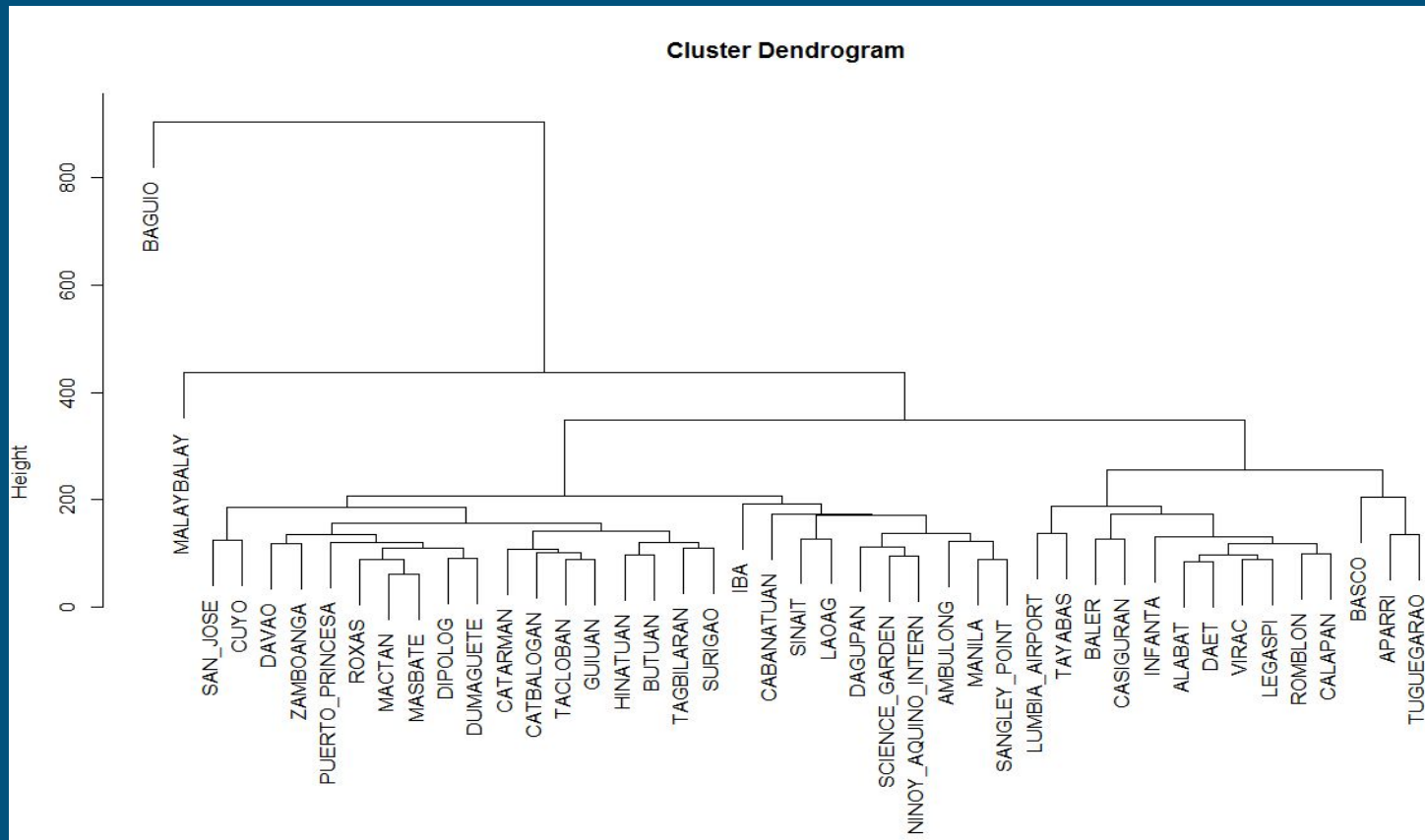
Distance matrix (dist) + hierarchical clustering (hclust) + split to subtrees (cutree)
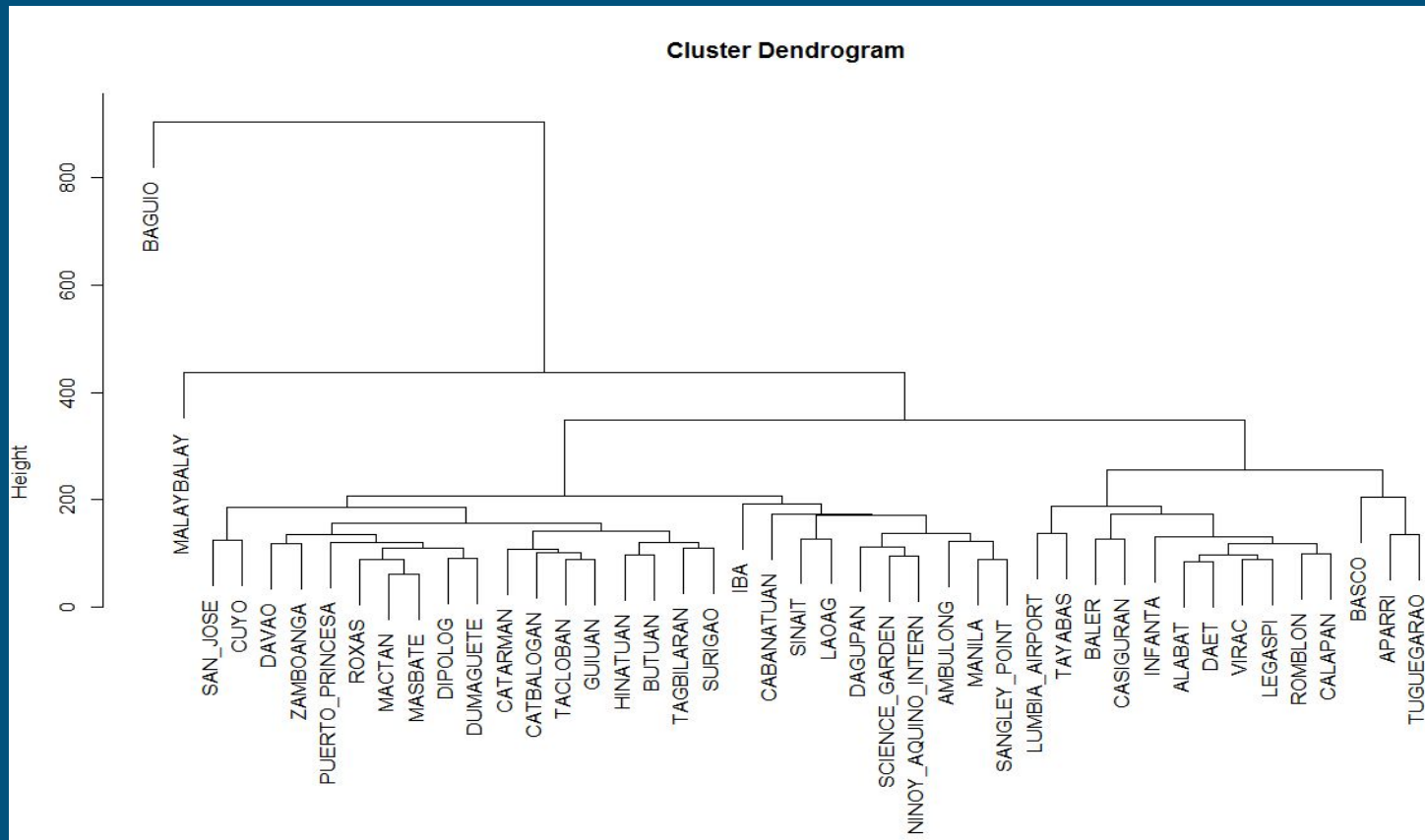
Temperature, Precipitation and Temperature + Precipitation

# Temperature



**Cluster Dendrogram**

# Some observations…

- Baguio is furthest (Summer Capital)
- Cooler by 8 °C (19 °C)

# Temperature



Cluster Dendrogram
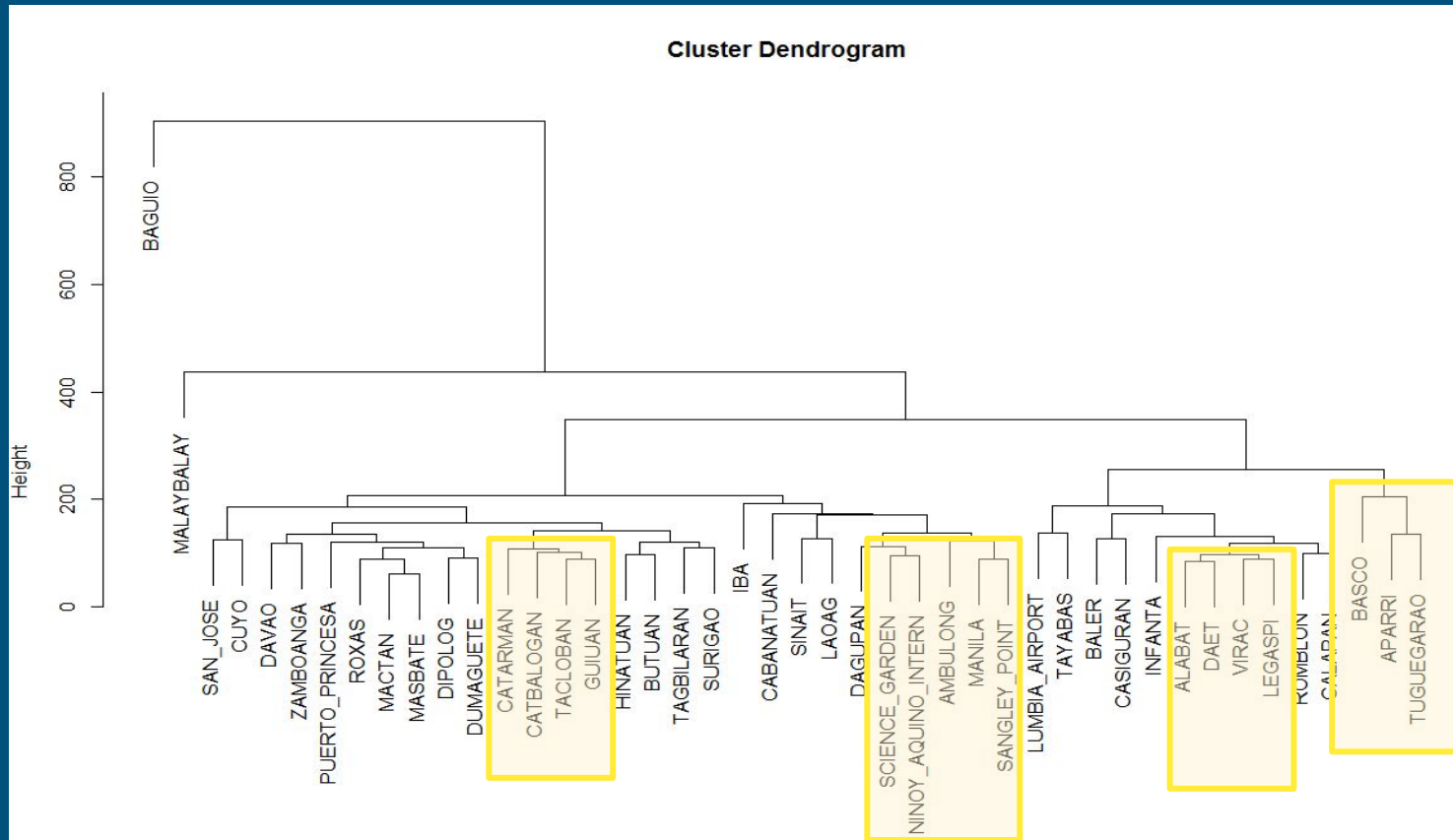
# Some observations…

- **WHY, Malaybalay?? (it rhymes)**
  - 2nd to Baguio at 24 °C

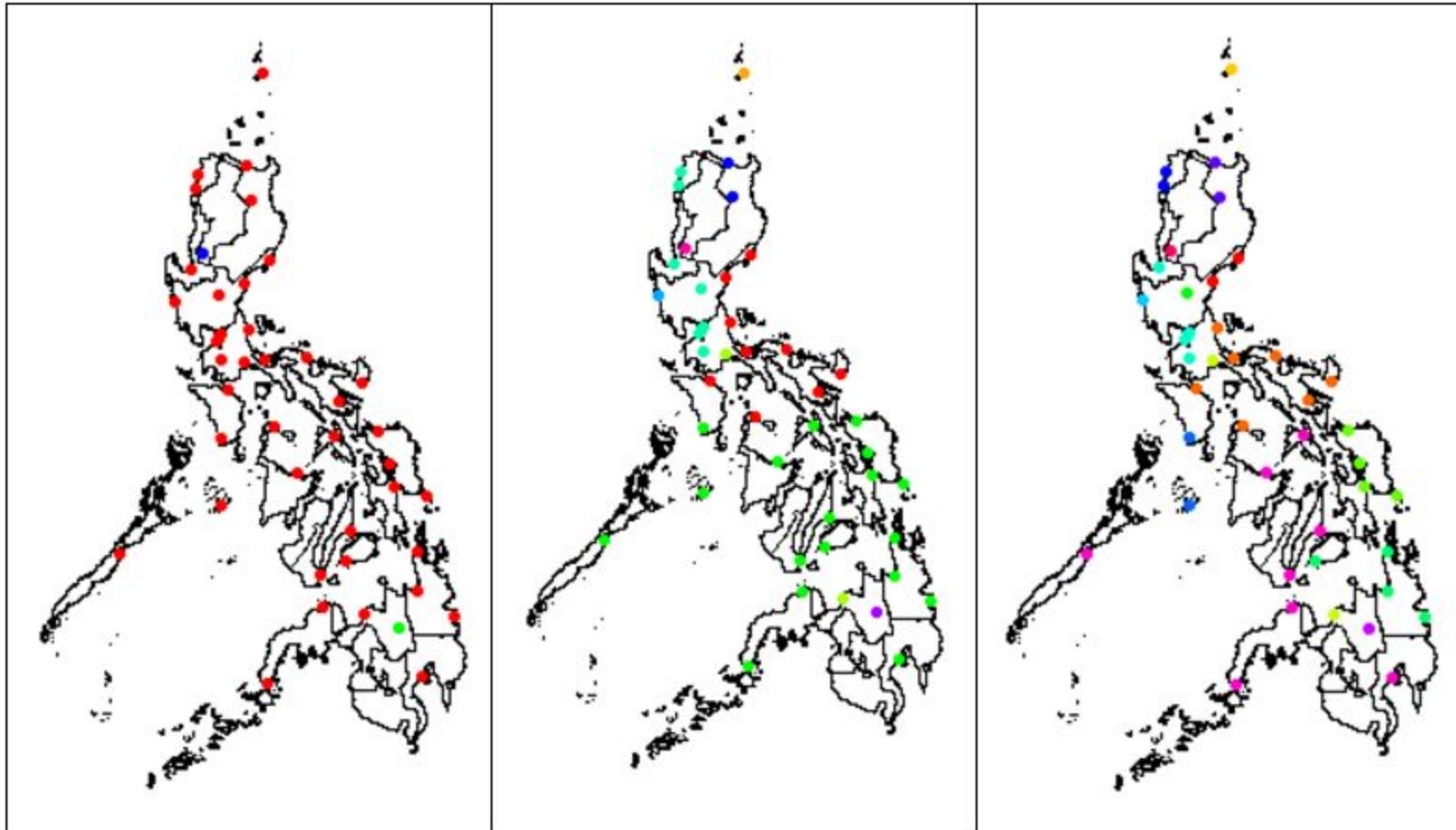"pleasant due to its altitude and the usual extreme heat of the tropical region is lacking"*

*http://www.bukidnon.gov.ph/home/index.php/about-bukidnon/general-info/climate
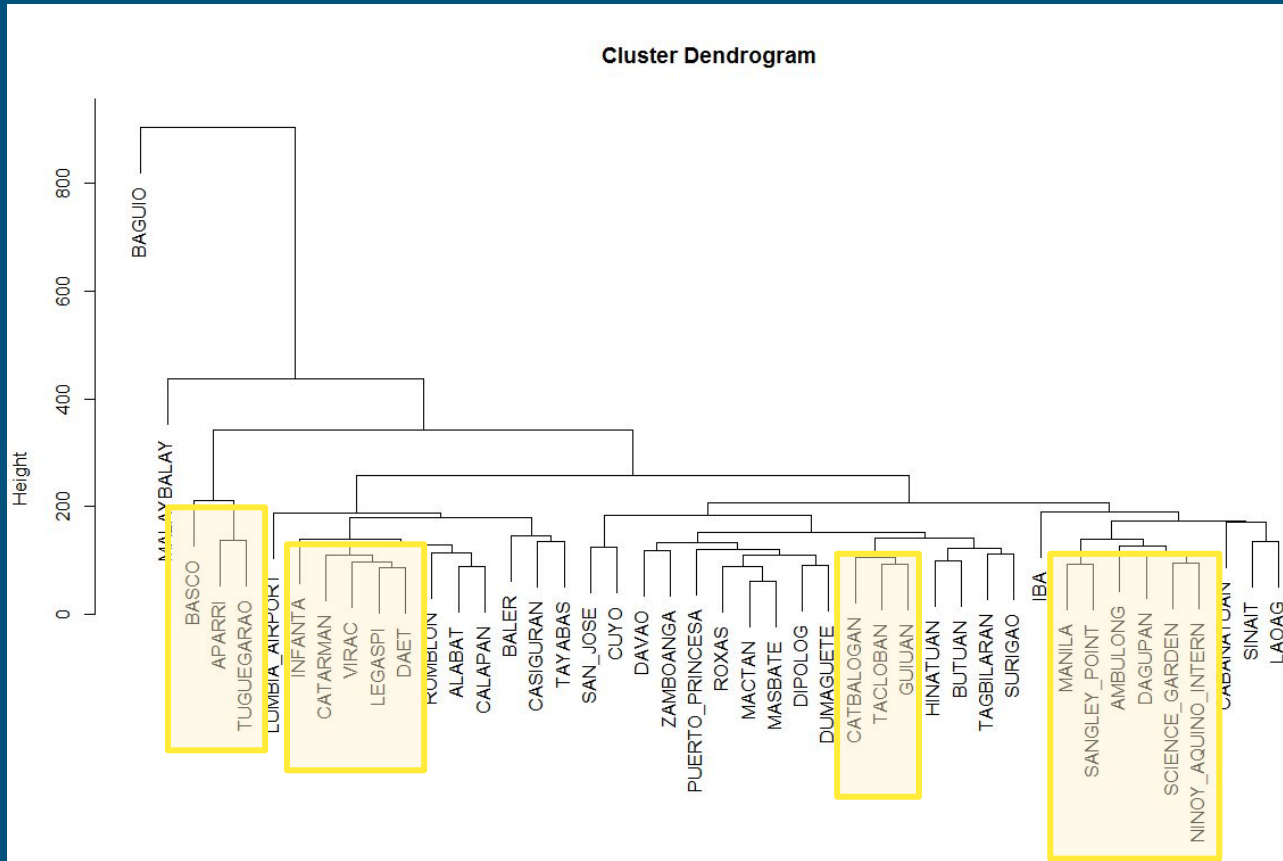
# Temperature



Cluster Dendrogram

# cutree

3, 9, 15

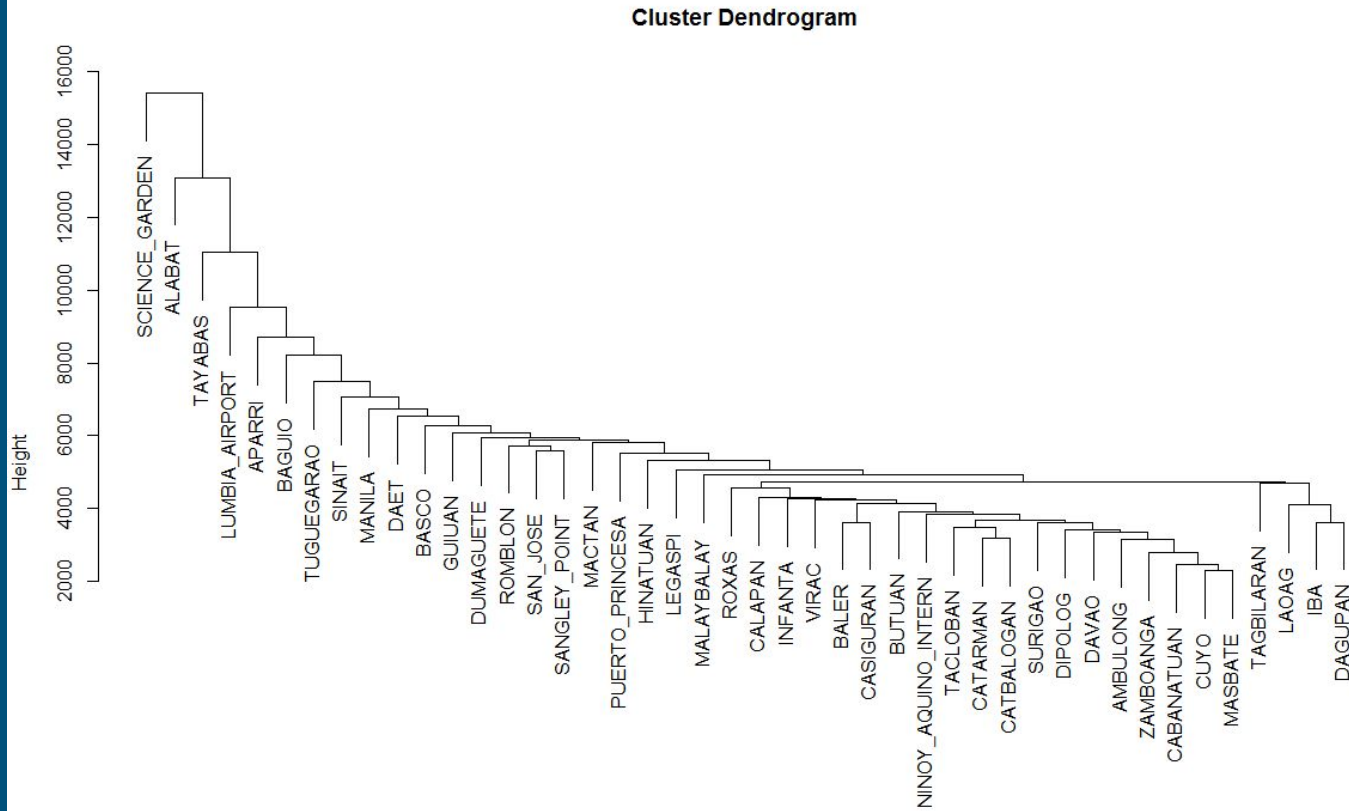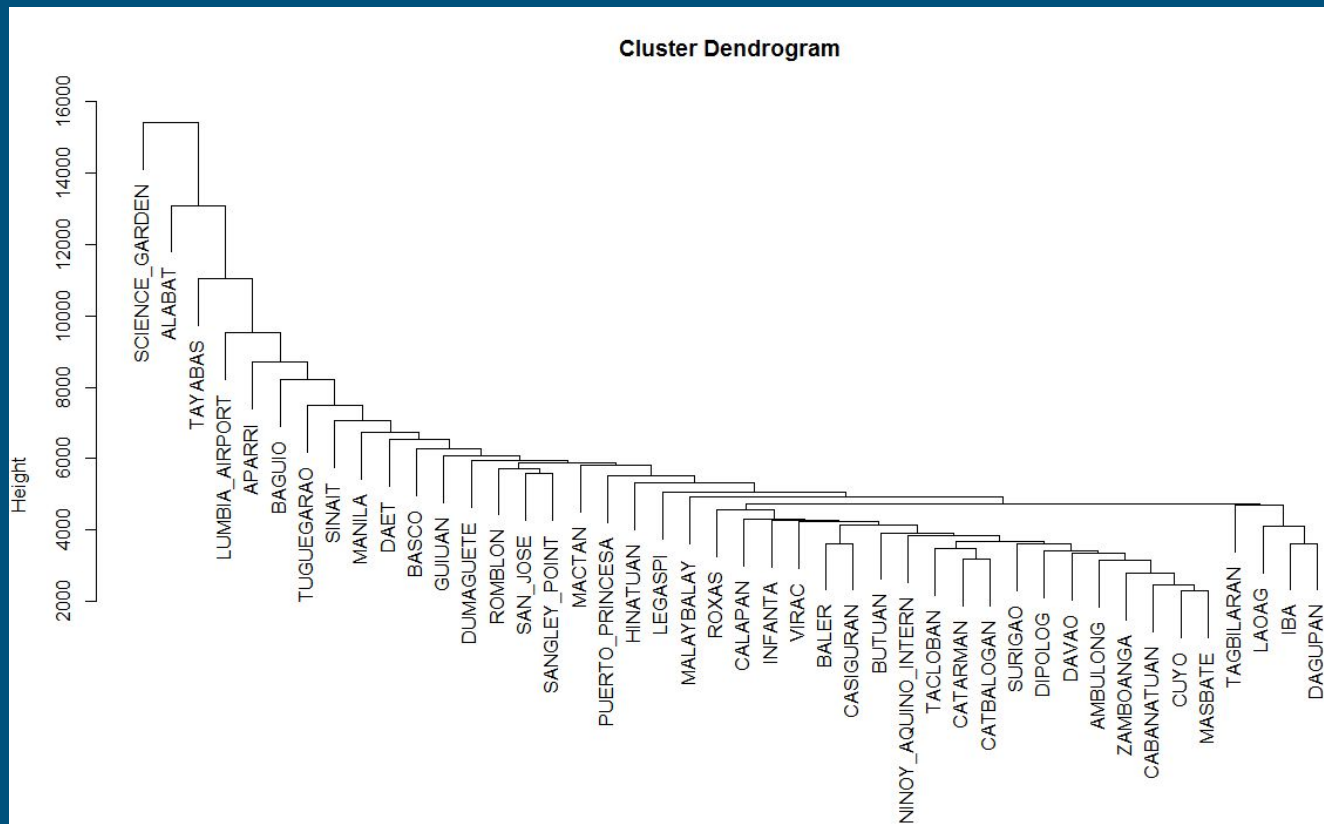# Temperature with no NAs

Replaced the NAs with the mean



Cluster Dendrogram

# cutree

3, 9, 15

# Precipitation



Cluster Dendrogram

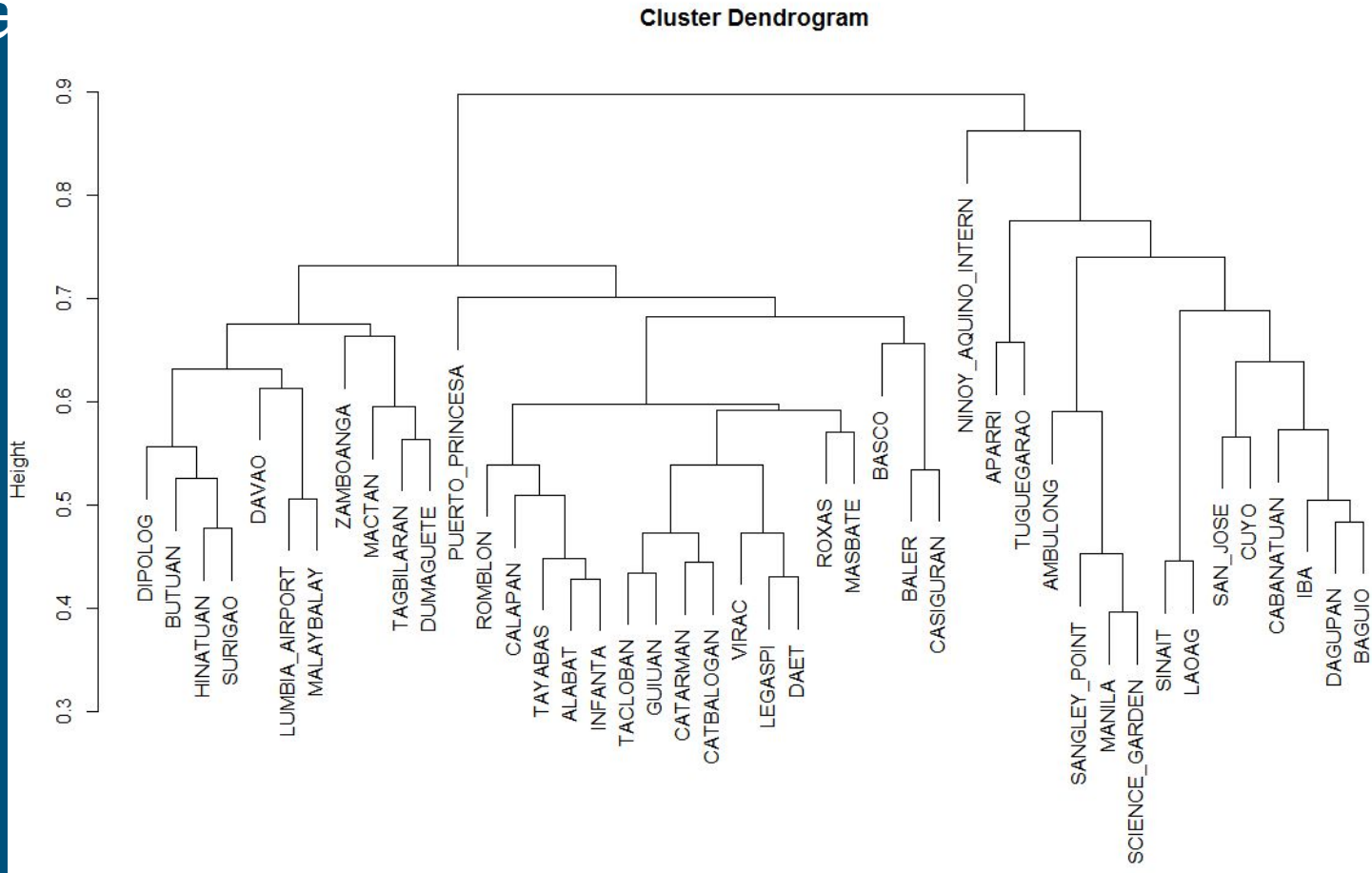# Precipitation with no NAs



Cluster Dendrogram

# Precipitation

Rain or no rain

Binary distance matrix

Rain (precipitation > 0) or no rain (precipitation = 0)
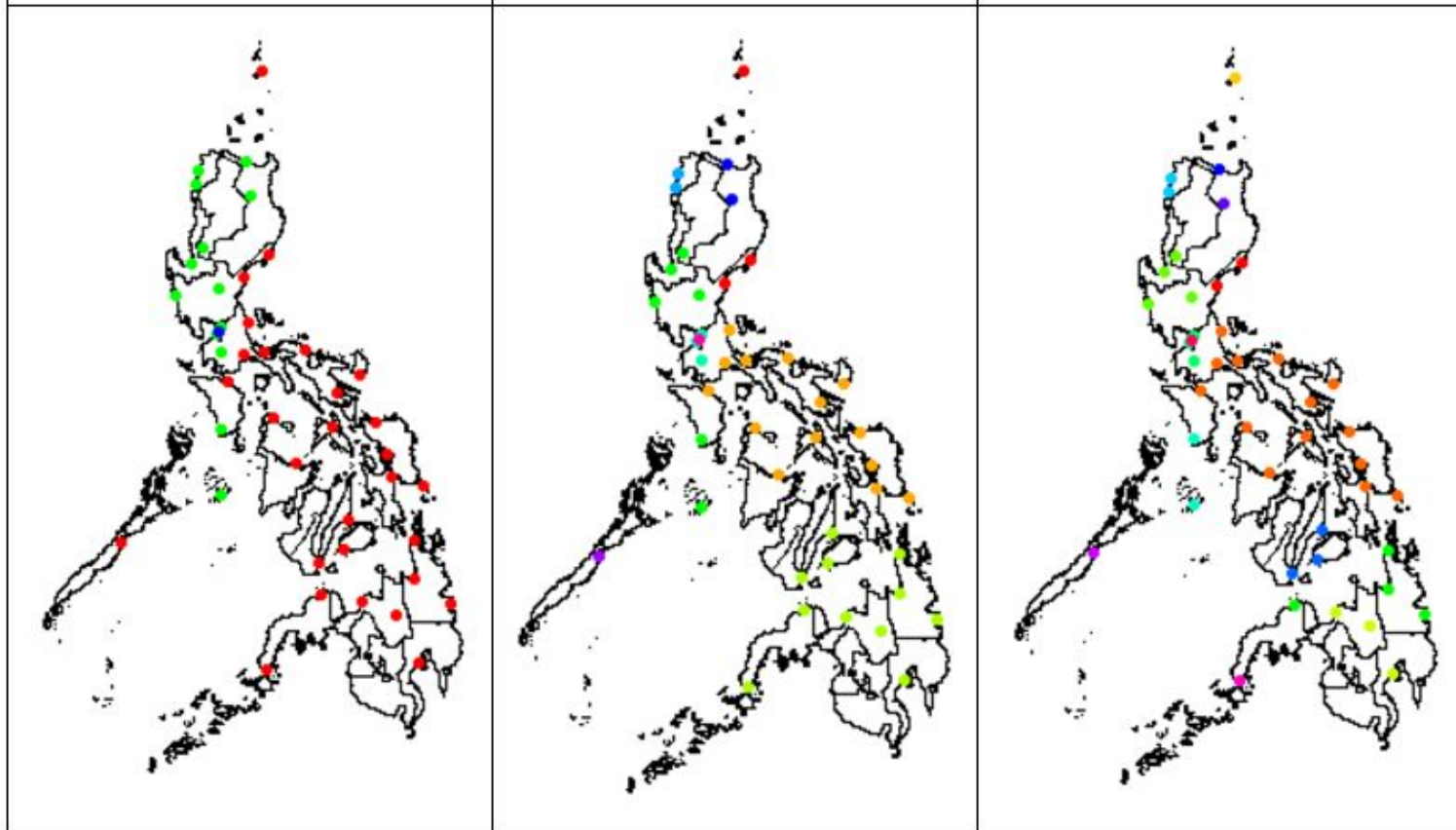
method="binary" in the dist function

# Precipita
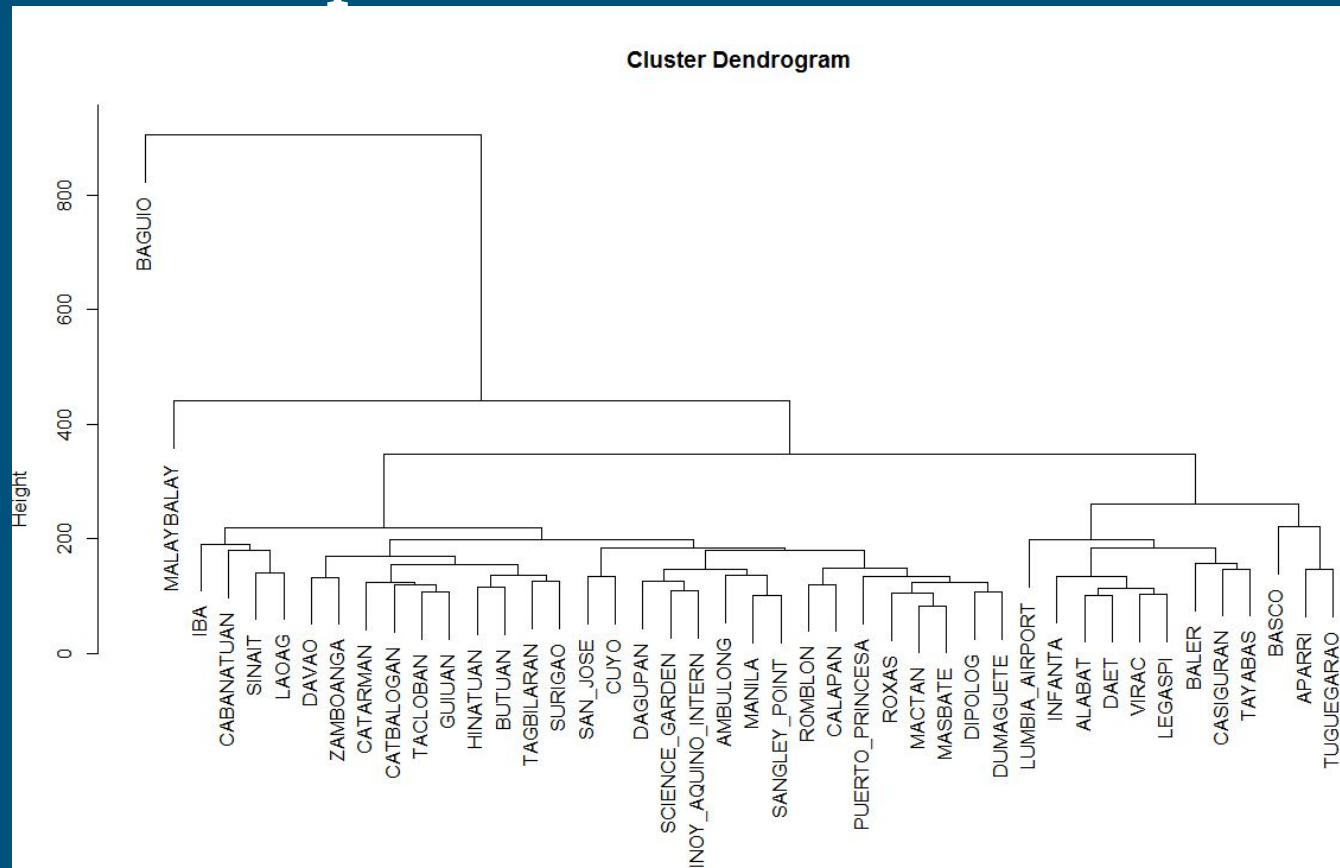


Cluster Dendrogram

# Precipitation

3, 9, 15

# Temperature + Precipitation

Converted the precipitation to binary

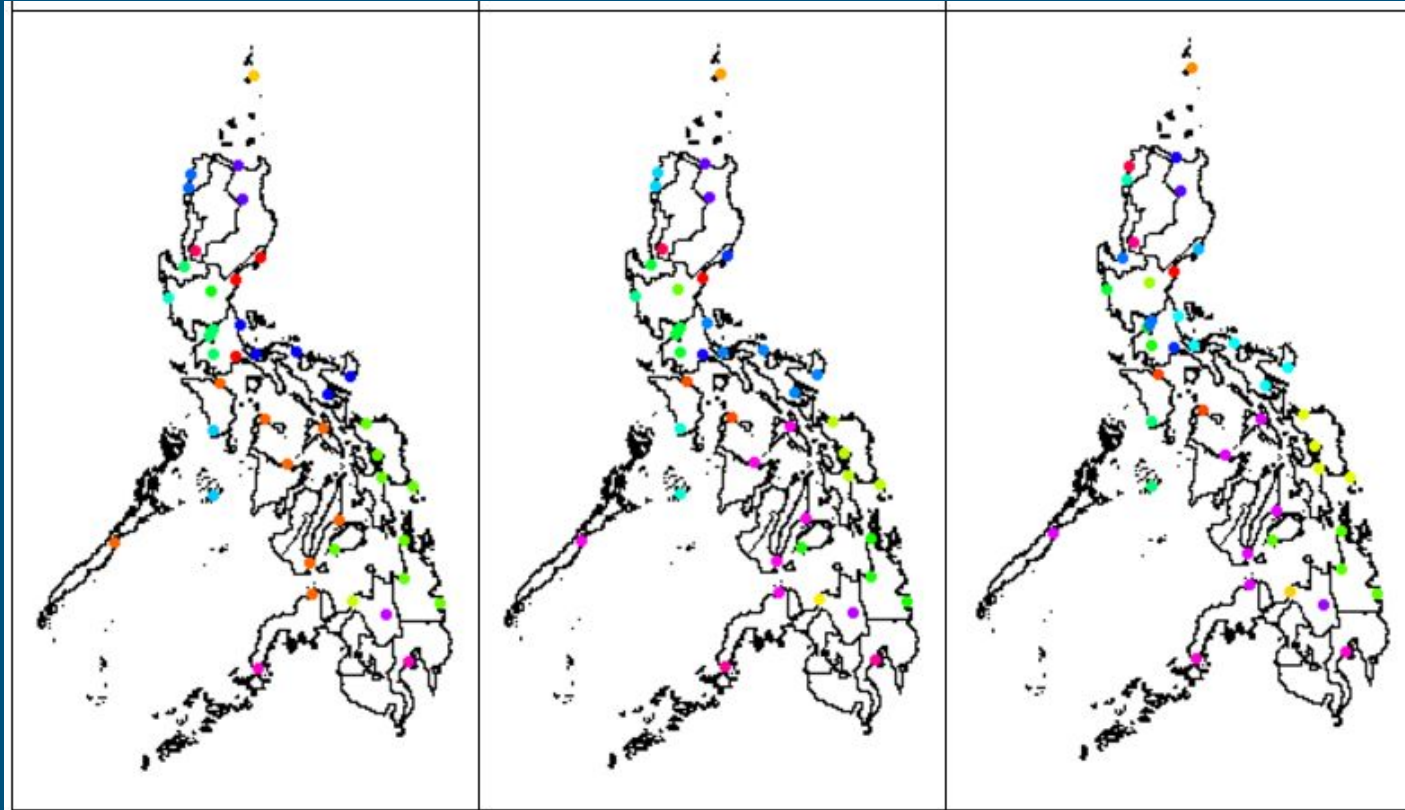Hclust + dist (complete distance measure) + cutree

# Temperature + Precipitation



**Cluster Dendrogram**

# Temperature + Precipitation
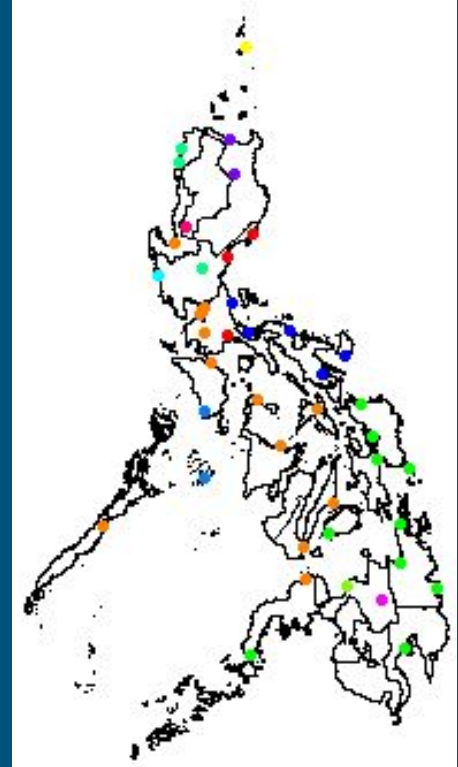
3, 7, 12

# Temperature + Precipitation

15, 19, 22

# Some observations...

- Eastern coast is "further" from the western coast
  - Bicol region and Quezon province
  - Samar and Leyte

Thank you