

# Weather Patterns X COVID-19

## Final Project Documentation

Uhuru Kamau, Shuyang Lu, Yifan Zhao, Chenhao Zhao, Zhaofeng Liu

## Contents

<b>Project Workflow</b>	<b>2</b>
Right Side . . . . .	2
Left Side . . . . .	3
<b>Data Acquisition</b>	<b>3</b>
1. New York City COVID-19 Data Archive . . . . .	3
2. New York City Weather Data . . . . .	4
3. Daily UV Index Scores - New York City . . . . .	4
<b>Relational Schema</b>	<b>5</b>
<b>Data Cleaning</b>	<b>6</b>
1. New York City COVID-19 Data Archive . . . . .	6
2. New York City Weather Data . . . . .	6
3. Daily UV Index Scores - New York City . . . . .	9

```
# Load Required Packages
library(tidyverse)
library(kableExtra)
library(readr)
library(gridExtra)
library(knitr)
```

# Project Workflow

## Right Side

```
knitr::include_graphics(path = "images/QBS181_1.png")
```



Figure 1: QBS181 Final Project: WorkFlow (Part A)

## Left Side

```
knitr::include_graphics(path = "images/QBS181_2.png")
```

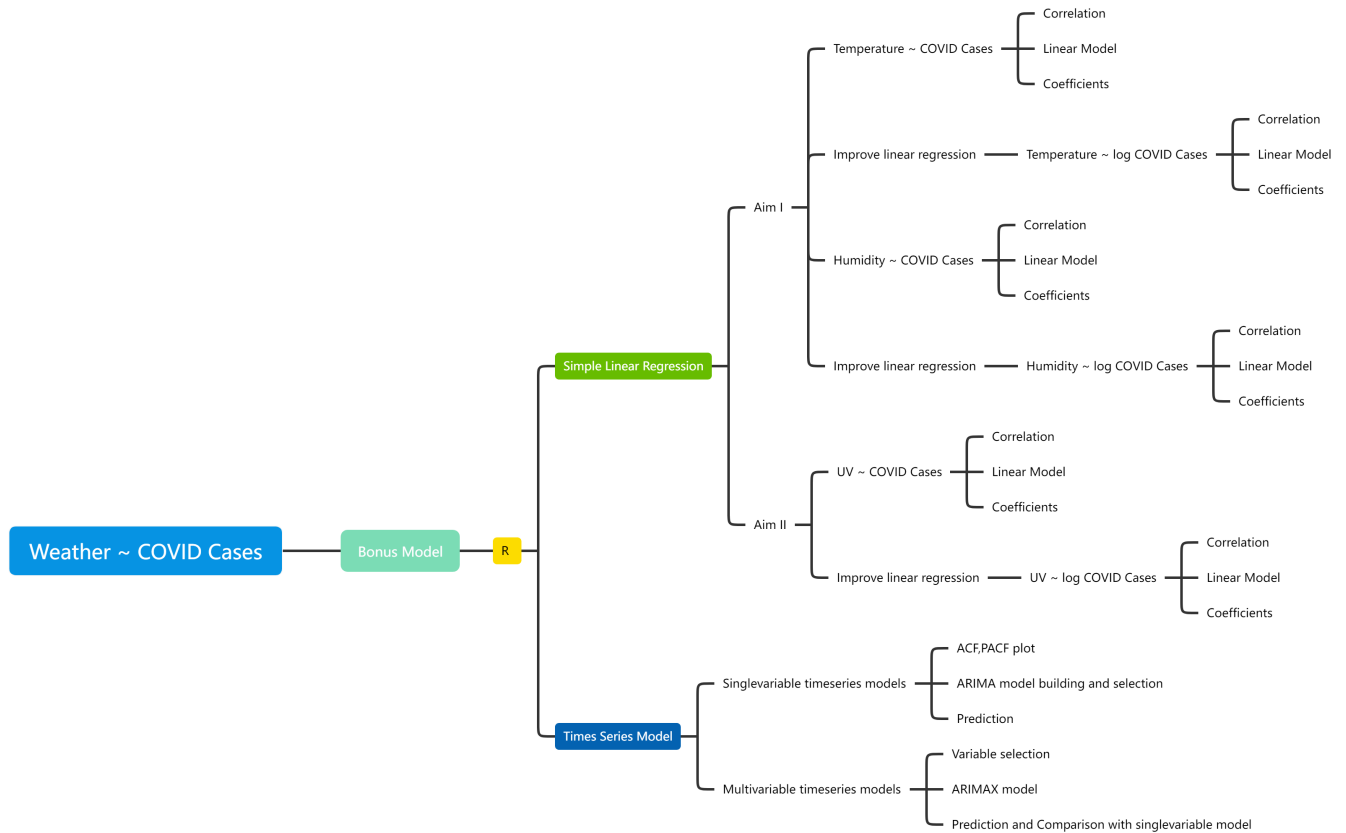


Figure 2: QBS181 Final Project: WorkFlow (Part B)

## Data Acquisition

### 1. New York City COVID-19 Data Archive

- Source: [NYC OpenData](#)
- Acquisition Method
  - Download .csv file
- Purpose:
  - We will use this time series data to track changes in the incidence of COVID-19.

## 2. New York City Weather Data

- Source: [Weather Underground - Weather Archive](#)
- Acquisition Method
  - Webscraping/ API Tool
- Purpose:
  - Merge time series weather data with timeseries Covid-19 data and investigate potential associations

## 3. Daily UV Index Scores - New York City

- Source: [Central New York's Live Weather Source](#)
- Acquisition Method
  - UV index values are presented as tables (see figure)
  - Copy tables and paste into Microsoft Excel
  - Save as .csv file
- Purpose
  - Sunlight and Vitamin-D absorbtion
    - \* It is generally accepted that there is a positive association between exposure to sunlight and absorbtion of vitamin-D.
    - \* It is also generally accepted that there is a positive association between vitamin-D absorbtion and immune system capacity.
  - We will us UV-Index as a proxy for exposure to sunlight at the population level and test for associations between UV Index and the incidence of Covid-19.

## Relational Schema

```
knitr::include_graphics(path = "images/Relational_Schema.png")
```

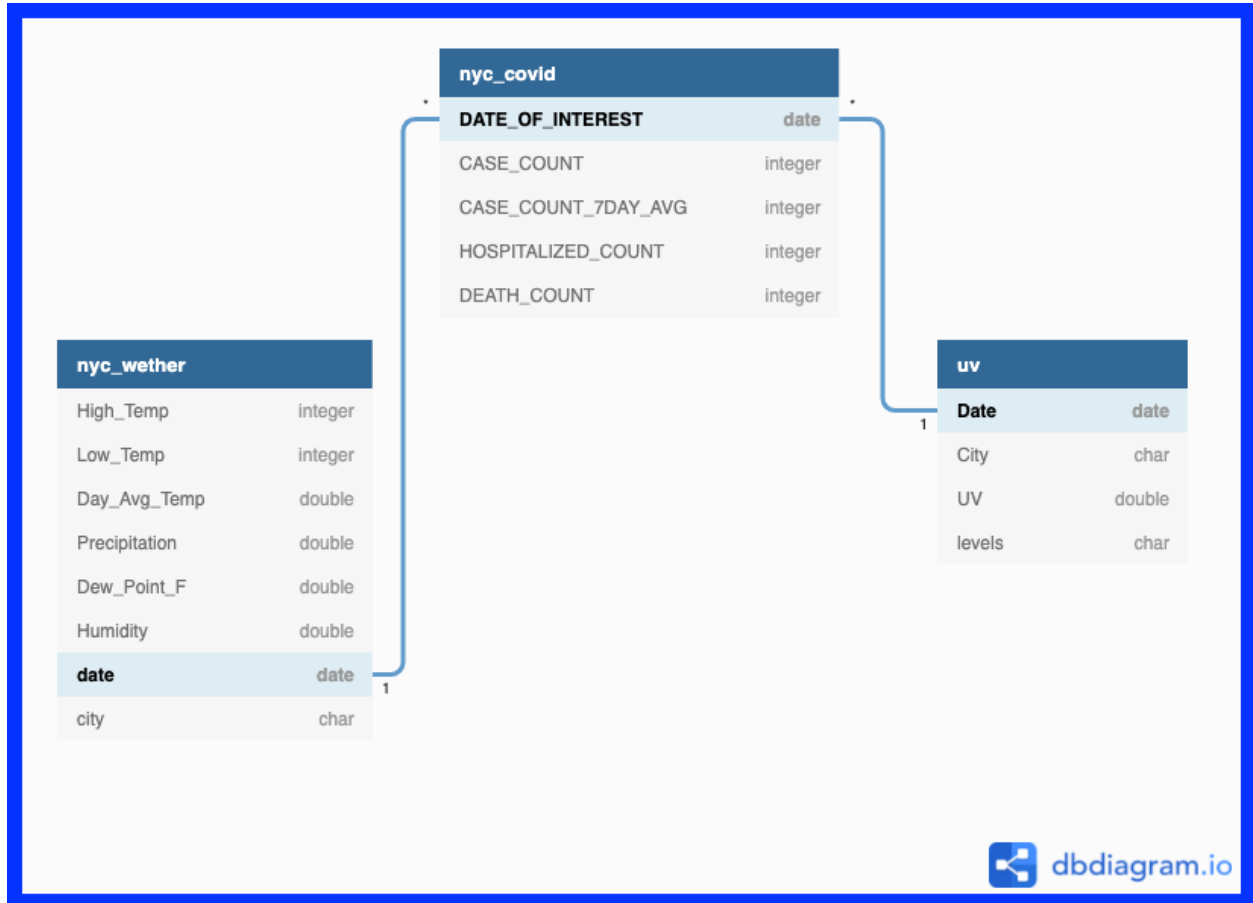


Figure 3: Highlighting the Keys to our Relational Database

# Data Cleaning

## 1. New York City COVID-19 Data Archive

### A. Read in the Covid-19 Date Frame

```
covid_df <- read.csv("data/Raw Data/nyc_covid19_data/NYC_Covid_Data_raw.csv")
```

## 2. New York City Weather Data

### A. Read-in File from Raw Data File

- The raw file has an issue with the column headers.
  - Several Headers include symbols that don't work with the interpreter
    - \* eg. *Low\_Temp( $^{\circ}F$ )*, *High\_Temp( $^{\circ}F$ )*
- **Solution:** Update column names while reading in the file!

```
# vector with acceptable column names
headers <- c("High.Temp", "Low.Temp", "Avg.Temp", "Precip", "Dew.Point", "Humidity",
            "Date", "City")

# read in the raw data file
weather.raw <- read.csv(file = "data/Raw Data/nyc_weather_raw.csv", header = TRUE,
                        col.names = headers)
```

### B. Format all observations of the “Date” Variable

- Variable is of class “character” by default

```
class(weather.raw$Date)
```

```
RESULT [1] "character"
```

- Reclassify the variable as a “Date”

```
weather.clean <- weather.raw %>%
  mutate(Date = as.Date(Date, "%m/%d/%Y"))
```

- Outcome:

```
class(weather.clean$Date)
```

```
RESULT [1] "Date"
```

### C. Missing Data

- Since we intend to do a time series, we need to identify any missing dates in the “date” column.
  - We will do this using a **CUSTOM FUNCTION!**

```
# Custom function to find the missing date in the date column
find.missing.dates <- function(d) {
  date_range <- seq(min(d), max(d), by = 1)
  date_range[!date_range %in% d]
}
```

- Use the custom function to identify missing dates in our NYC Weather df

```
# Display the missing dates
date.missing = c()
date.missing <- find.missing.dates(weather.clean$Date)
print(date.missing)
```

RESULT [1] "2020-11-08"

#### D. Replace Missing Values

- Method: Fill the missing data by averaging the former 6 days' data

```
# Find the index of the day before '2020-11-08'
weather.clean$Date <- as.character(weather.clean$Date)
id.missing.date = which(weather.clean$Date == "2020-11-07") + 1
```

- Build a custom function to fill the missing data
  - Approach: use the average of the previous six days

```
# Custom function to fill the missing data by averaging the former 6 days' data
fill.missing.values <- function(df, newrow.id) {
  newrow <- list()
  value <- c()
  first.row = newrow.id - 6
  last.row = newrow.id - 1
  col.num = ncol(df) - 2
  for (i in 1:col.num) {
    subs <- weather.clean[first.row:last.row, i] # Create a new subset for each column
    value <- mean(subs) # Calculate the mean
    newrow <- append(newrow, value)
  }
  return(newrow)
}
```

- Use the custom function to fill in the values of the missing row

```
# Fill the missing values in the missing row
missing.row <- fill.missing.values(weather.clean, id.missing.date)
missing.row <- append(missing.row, "2020-11-08")
missing.row <- append(missing.row, "new york city")
```

- Build another custom function to insert the row into the df

```
# Custom function to insert the new row
insertRow <- function(existingDF, newrow, r) {
  existingDF[seq(r + 1, nrow(existingDF) + 1), ] <- existingDF[seq(r, nrow(existingDF)),
    ]
  existingDF[r, ] <- newrow
  existingDF
}
```

- Insert the imputed value into the df!

```
# Insert the missing row and store it into a new df
weather.clean <- insertRow(weather.clean, missing.row, id.missing.date)
```

#### E. Remove the “City” Variable

- Every observation is is “new york city”
  - This variable is effectively just clutter.

```
weather.clean = weather.raw %>%
  select(-City)
```

#### F. Display

```
kable(x = weather.clean[1:5, ], digits = 2, align = "c")
```

High.Temp	Low.Temp	Avg.Temp	Precip	Dew.Point	Humidity	Date
44	26	35.46	0.00	13.67	41.83	3/1/2020
56	38	48.17	0.00	30.46	51.12	3/2/2020
58	48	52.41	0.01	44.59	75.47	3/3/2020
57	46	50.52	0.28	28.52	44.76	3/4/2020
52	40	44.75	0.00	25.38	48.50	3/5/2020

#### G. Write the Processed Data to a new .csv file

```
write.csv(x = weather.clean, file = "data/Processed Data/nyc_weather.csv")
```

#### H. In Excel - Add a “Season” variable to the .csv generated in the prior section

- Open File in Microsoft Excel
  - Summarise the process
- Save the file as “nyc\_clean\_weather\_add\_season.csv”
  - Push to Github

#### I. Read in “nyc\_clean\_weather\_add\_season.csv” and reformat the “date” variable

- Read in the file



```
# vector with acceptable column names
headers <- c("High.Temp", "Low.Temp", "Avg.Temp", "Precip", "Dew.Point", "Humidity",
  "date", "City", "Month", "Season")
add_season_nyc_weather <- read.csv(file = "data/nyc_clean_weather_add_season.csv",
  header = T, col.names = headers) %>%
  select(-City)
```

- Reformat the “date” variable

```
# change date type from ymd to mdy
add_season_nyc_weather$date <- format(as.Date(add_season_nyc_weather$date, "%Y/%m/%d"),
  "%m/%d/%Y")
```

J. Display

```
knitr::kable(x = add_season_nyc_weather[1:5, ], align = "c")
```

High.Temp	Low.Temp	Avg.Temp	Precip	Dew.Point	Humidity	date	Month	Season
44	26	35.46	0.00	13.67	41.83333	03/01/2020	3	spring
56	38	48.17	0.00	30.46	51.12500	03/02/2020	3	spring
58	48	52.41	0.01	44.59	75.47059	03/03/2020	3	spring
57	46	50.52	0.28	28.52	44.76000	03/04/2020	3	spring
52	40	44.75	0.00	25.38	48.50000	03/05/2020	3	spring

### 3. Daily UV Index Scores - New York City

A. Read in the csv file that we built in excel

```
nyc.uv = read.csv(file = "data/Raw Data/nyc_uv_raw.csv", header = T)
```

B. Remove the “City” Variable

- Every observation is “New York City” so this variable is effectively clutter.

```
nyc.uv <- nyc.uv %>%
  select(-City)
```

C. Display

```
knitr::kable(x = nyc.uv[1:5, ], align = "c")
```

Date	UV
3/1/2020	0.8
3/2/2020	1.3
3/3/2020	1.3
3/4/2020	0.8
3/5/2020	2.1