

# **Is Rainfall Getting Heavier?**

## **Building a Weather Forecasting Pipeline with Singapore Weather Station Data**

By: Chin Hwee Ong (@ongchinhwee)

7 February 2021  
FOSDEM Python devroom

# About me

Ong Chin Hwee 王敬惠

- Data Engineer
- Based in **sunny Singapore** 🌞
- Aerospace Engineering +  
Computational Modelling
- Loves (and contributes to) pandas



@ongchinhwee

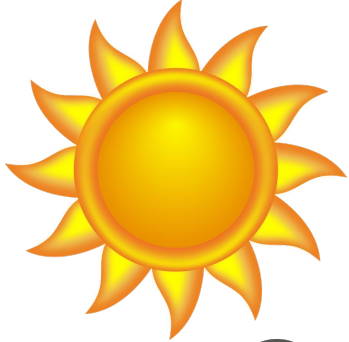
Singapore 新加坡:  
1°17'22.81"N, 103°  
51'0.25"E  
北纬1度, 经纬103度





Singapore is a **tropical**  
country

We have our “four seasons”:



1. **Cold and Rainy**

2. Warm and Dry

3. **Extremely Hot**

4. **Hot and Stormy**



# PONDING

BECAUSE FLOODS ONLY HAPPEN ONCE EVERY 50 YEARS

**Since 2018, Singapore had more than 20 flash floods.**

**Majority of the floods were caused by intense rain.**

Source: PUB Singapore

(<https://www.pub.gov.sg/drainage/floodmanagement/recentflashfloods>)

@ongchinwee



Could we **predict heavier rainfall**  
with weather data?

# Extracting Weather Data



## Realtime Weather Readings across Singapore

### FILES IN THIS DATASET

#### Air Temperature across Singapore



Views:



< > Embed Chart

API View

#### Rainfall across Singapore



#### Relative Humidity across Singapore



#### Wind Direction across Singapore



GET

<https://api.data.gov.sg/v1/environment/air-temperature>

Get air temperature readings across Singapore

- Has per-minute readings from NEA
- Use the `date_time` parameter to retrieve the latest available data at that moment in time
- Use the `date` parameter to retrieve all of the readings for that day.

Parameters

Try it out

Name

Description

People

Data.gov.sg - Singapore's Open Data Portal

@ongchinwee

# Realtime Weather Readings across Singapore

**Real-time API** on Data.gov.sg (Singapore's open data portal)

**Open government data** available under the Singapore Open Data License

**(Almost) minute-by-minute** weather station readings

“Let’s try to scrap weather data for a specific weather station!”

“How about we scrap multi-day data  
from the API?”

date\_time

YYYY-MM-DD[T]HH:mm:ss (SGT)

string

(query)

date\_time - YYYY-MM-DD[T]HH:mm:ss (SGT)

date

YYYY-MM-DD

string

(query)

2019-11-23

Execute

Clear

Responses

Curl

```
curl -X GET "https://api.data.gov.sg/v1/environment/air-temperature?date=2019-11-23" -H "accept: application/json"
```

Request URL

```
https://api.data.gov.sg/v1/environment/air-temperature?date=2019-11-23
```

## Responses

### Curl

```
curl -X GET "https://api.data.gov.sg/v1/environment/air-temperature?date=2019-11-23" -H "accept: application/json"
```

### Request URL

```
https://api.data.gov.sg/v1/environment/air-temperature?date=2019-11-23
```

```
    "id": "S100",
    "device_id": "S100",
    "name": "Woodlands Road",
    "location": {
      "latitude": 1.4172,
      "longitude": 103.74855
    }
  },
  {
    "id": "S115",
    "device_id": "S115",
    "name": "Tuas South Avenue 3",
    "location": {
      "latitude": 1.29377,
      "longitude": 103.61843
    }
  }
],
"reading_type": "DBT 1M F",
"reading_unit": "deg C"
},
"items": [
  {
    "timestamp": "2019-11-23T00:01:00+08:00",
    "readings": [
      {
        "station_id": "S109",
        "value": 25.2
      },
      {
        "station_id": "S117",
        "value": 26
      }
    ]
  }
],
```



```
    "id": "S100",
    "device_id": "S100",
    "name": "Woodlands Road",
    "location": {
      "latitude": 1.4172,
      "longitude": 103.74855
    }
  },
  {
    "id": "S115",
    "device_id": "S115",
    "name": "Tuas South Avenue 3",
    "location": {
      "latitude": 1.29377,
      "longitude": 103.61843
    }
  }
],
"reading_type": "DBT 1M F",
"reading_unit": "deg C"
},
"items": [
  {
    "timestamp": "2019-11-23T00:01:00+08:00",
    "readings": [
      {
        "station_id": "S109",
        "value": 25.2
      },
      {
        "station_id": "S117",
        "value": 26
      }
    ]
  }
],
```

# Nested JSON format!

@ongchinwee

```

    "id": "S100",
    "device_id": "S100",
    "name": "Woodlands Road",
    "location": {
      "latitude": 1.4172,
      "longitude": 103.74855
    }
  },
  {
    "id": "S115",
    "device_id": "S115",
    "name": "Tuas South Avenue 3",
    "location": {
      "latitude": 1.29377,
      "longitude": 103.61843
    }
  }
],
"reading_type": "DBT 1M F",
"reading_unit": "deg C"
},
"items": [
  {
    "timestamp": "2019-11-23T00:01:00+08:00",
    "readings": [
      {
        "station_id": "S109",
        "value": 25.2
      },
      {
        "station_id": "S117",
        "value": 26
      }
    ]
  }
]

```



|     | index                  | readings  | timestamp                 |
|-----|------------------------|---|---------------------------|
| [ ] |                        |   |                           |
| ↳   | 0                      | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:05:00+08:00 |
|     | 1                      | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:10:00+08:00 |
|     | 2                      | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:15:00+08:00 |
|     | 3                      | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-06-01 00:20:00+08:00 |
|     | 4                      | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:25:00+08:00 |
|     | ...                    | ...   | ...                       |
|     | 60435                  | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-12-31 23:35:00+08:00 |
|     | 60436                  | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:40:00+08:00 |
|     | 60437                  | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-12-31 23:45:00+08:00 |
|     | 60438                  | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:50:00+08:00 |
|     | 60439                  | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:55:00+08:00 |
|     | 60440 rows x 3 columns |   |                           |

hweecat / api\_scraping\_nea\_datasets

Watch

1

★ Star

1

🔗 Fork

1

<> Code

🔔 Issues 0

🔗 Pull requests 0

📁 Projects 0

🛡 Security

📊 Insights

## Join GitHub today

Dismiss

GitHub is home to over 40 million developers working together to host and review code, manage projects, and build software together.

Sign up

### Scraping Meteorological Data from Data.gov.sg APIs

📁 24 commits

🔗 2 branches

📦 0 packages

📦 0 releases

👤 2 contributors

Branch: master ▼

New pull request

Find file

Clone or download ▼



hweecat Merge pull request #8 from hweecat/airtemp\_rain ...

Latest commit 823bc3b 5 days ago



.gitignore

update timezone code for pandas 0.25

22 days ago



API\_scraping\_datagovsg\_(airtemp\_rainfall).py

add try-except logic for null-data-for-date case

7 days ago

## “Scraping Meteorological Data from Data.gov.sg APIs” Project

@ongchinhwee

# Data.gov.sg Weather Data API Scraping

Scraping weather data from APIs via “Requests” library

“Requests”:

Python library for humans to send HTTP requests



# Data.gov.sg Weather Data API Scraping

Currently supported Data.gov.sg APIs:

1. Air Temperature (in °C)
2. Rainfall (in mm)
3. Relative Humidity
4. Wind Direction
5. Wind Speed

Scrap data for continuous time range + specific weather station

# Design Considerations

## **Slow connection**

- retry mechanism

```
from retrying import retry
```

```
@retry(wait_exponential_multiplier=1000, wait_exponential_max=10000)
```

```
def get_rainfall_data_from_date(date):
```

# Design Considerations

## Slow connection

API working but **no data for specific date**

| Code | Details  |
|------|--|
| 200  | <p>Response body</p> <p>Download</p> <pre>{   "metadata": {     "stations": []   },   "items": [],   "api_info": {     "status": "healthy"   } }</pre> |

# Design Considerations

## Slow connection

API working but **no data for specific date**

- Return **empty DataFrame** with same column names as if there were data for specific date



# Design Considerations

## Slow connection

API working but **no data for specific date**

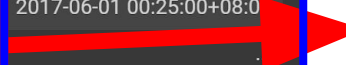
**Nested JSON** to pandas **DataFrame** conversion

- Extract desired **station** and **readings**
- Concatenate them back with **timestamp**

# Design Considerations

## Nested JSON to pandas DataFrame conversion

| index |     | readings  | timestamp                 |
|-------|-----|---|---------------------------|
| 0     | 0   | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:05:00+08:00 |
| 1     | 1   | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:10:00+08:00 |
| 2     | 2   | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:15:00+08:00 |
| 3     | 3   | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-06-01 00:20:00+08:00 |
| 4     | 4   | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:25:00+08:00 |
| ...   | ... | ...   | ...                       |
| 60435 | 277 | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-12-31 23:35:00+08:00 |
| 60436 | 278 | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:40:00+08:00 |
| 60437 | 279 | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-12-31 23:45:00+08:00 |
| 60438 | 280 | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:50:00+08:00 |
| 60439 | 281 | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:55:00+08:00 |



| station_id |      | value |
|------------|------|-------|
| 0          | S109 | 94.3  |
| 1          | S117 | 89.2  |
| 2          | S50  | 87.4  |
| 3          | S107 | 85.3  |
| 4          | S43  | 85.9  |
| ...        | ...  | ...   |
| 12         | S100 | 83.9  |
| 13         | S106 | NaN   |
| 14         | S115 | NaN   |
| 15         | S44  | NaN   |
| 16         | S122 | NaN   |

# Design Considerations

## Nested JSON to pandas DataFrame conversion

|       | index | readings  | timestamp                 |
|-------|-------|---|---------------------------|
|       | 0     | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:05:00+08:00 |
|       | 1     | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:10:00+08:00 |
|       | 2     | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:15:00+08:00 |
|       | 3     | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-06-01 00:20:00+08:00 |
|       | 4     | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-06-01 00:25:00+08:00 |
|       | ...   | ...   | ...                       |
| 60435 | 277   | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-12-31 23:35:00+08:00 |
| 60436 | 278   | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:40:00+08:00 |
| 60437 | 279   | [{'station_id': 'S77', 'value': 0}, {'station_... | 2017-12-31 23:45:00+08:00 |
| 60438 | 280   | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:50:00+08:00 |
| 60439 | 281   | [{'station_id': 'S77', 'value': 0.2}, {'statio... | 2017-12-31 23:55:00+08:00 |

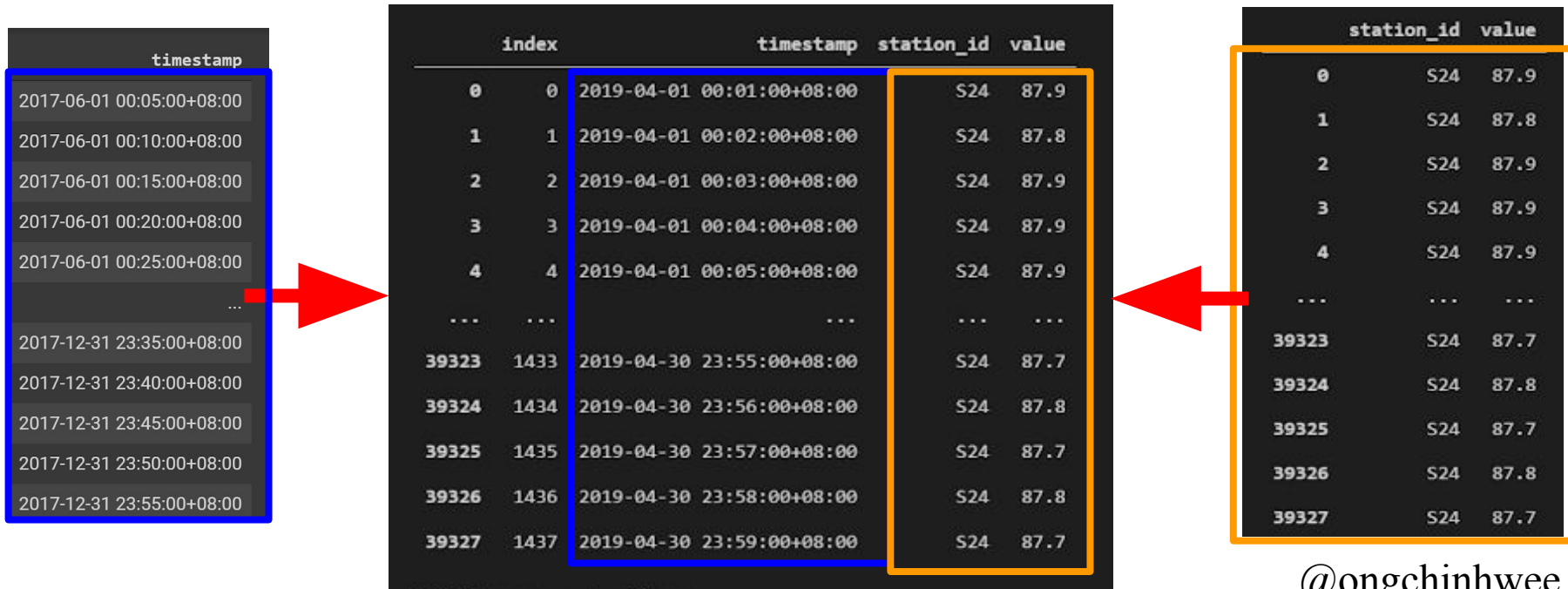
|     | station_id | value |
|-----|------------|-------|
| 0   | S109       | 94.3  |
| 1   | S117       | 89.2  |
| 2   | S50        | 87.4  |
| 3   | S107       | 85.3  |
| 4   | S43        | 85.9  |
| ... | ...        | ...   |
| 12  | S100       | 83.9  |
| 13  | S106       | NaN   |
| 14  | S115       | NaN   |
| 15  | S44        | NaN   |
| 16  | S122       | NaN   |



|       | station_id | value |
|-------|------------|-------|
| 0     | S24        | 87.9  |
| 1     | S24        | 87.8  |
| 2     | S24        | 87.9  |
| 3     | S24        | 87.9  |
| 4     | S24        | 87.9  |
| ...   | ...        | ...   |
| 39323 | S24        | 87.7  |
| 39324 | S24        | 87.8  |
| 39325 | S24        | 87.7  |
| 39326 | S24        | 87.8  |
| 39327 | S24        | 87.7  |

# Design Considerations

## Nested JSON to pandas DataFrame conversion



# **Singapore Rainfall Data: A 4-Year Time Series Analysis**

# Time Series Analysis of Singapore Rainfall Data



Selected **weather station**:

Changi Weather Station (**ID: S24**)

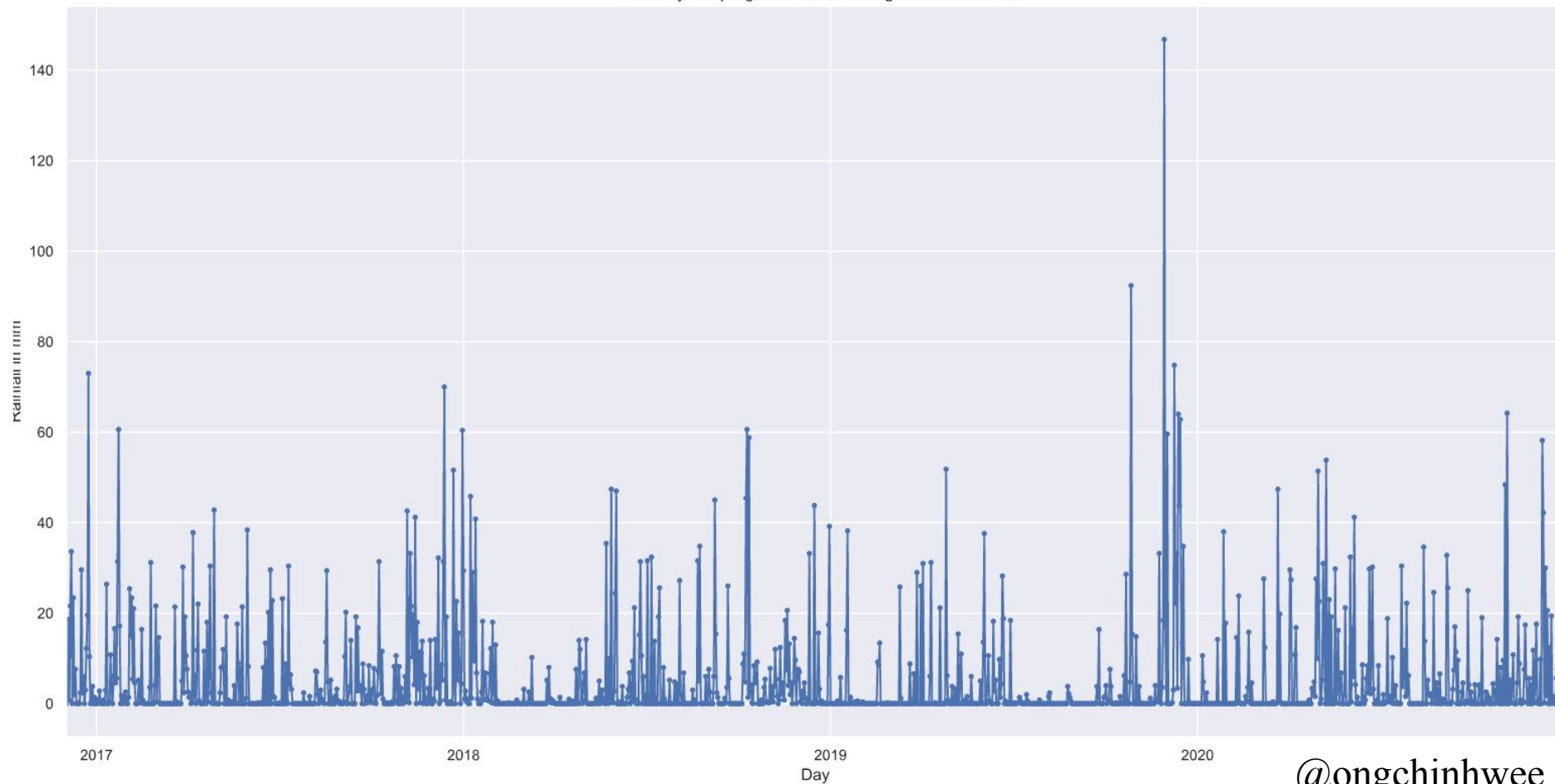
Analysis **timeframe**:

2 Dec 2016 to 31 Dec 2020 (**~4 years**)

**Objective:**

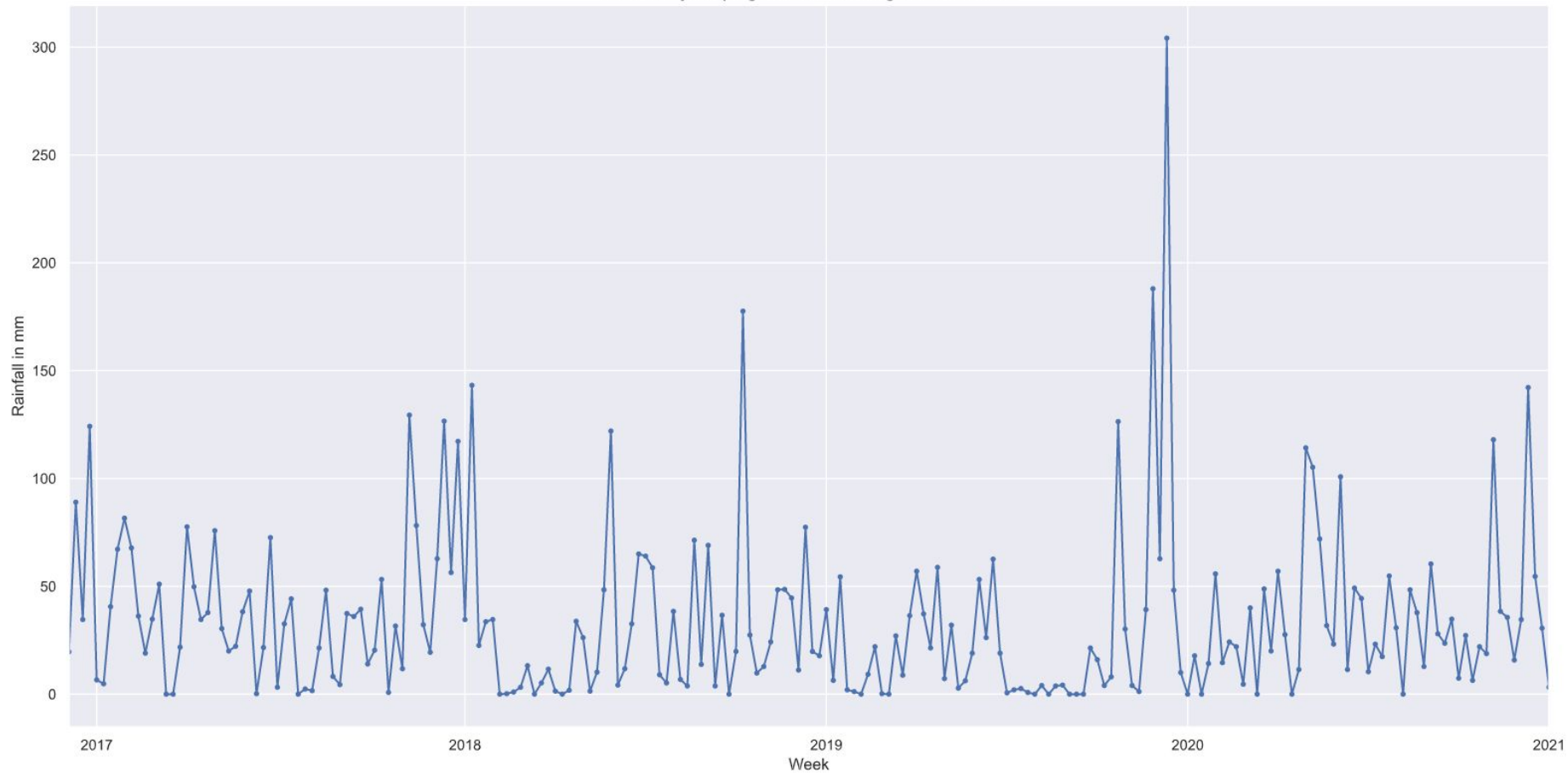
- Extract **trend** and **seasonality** from 5-minute rainfall data

Total daily sampling of rainfall at Changi Weather Station



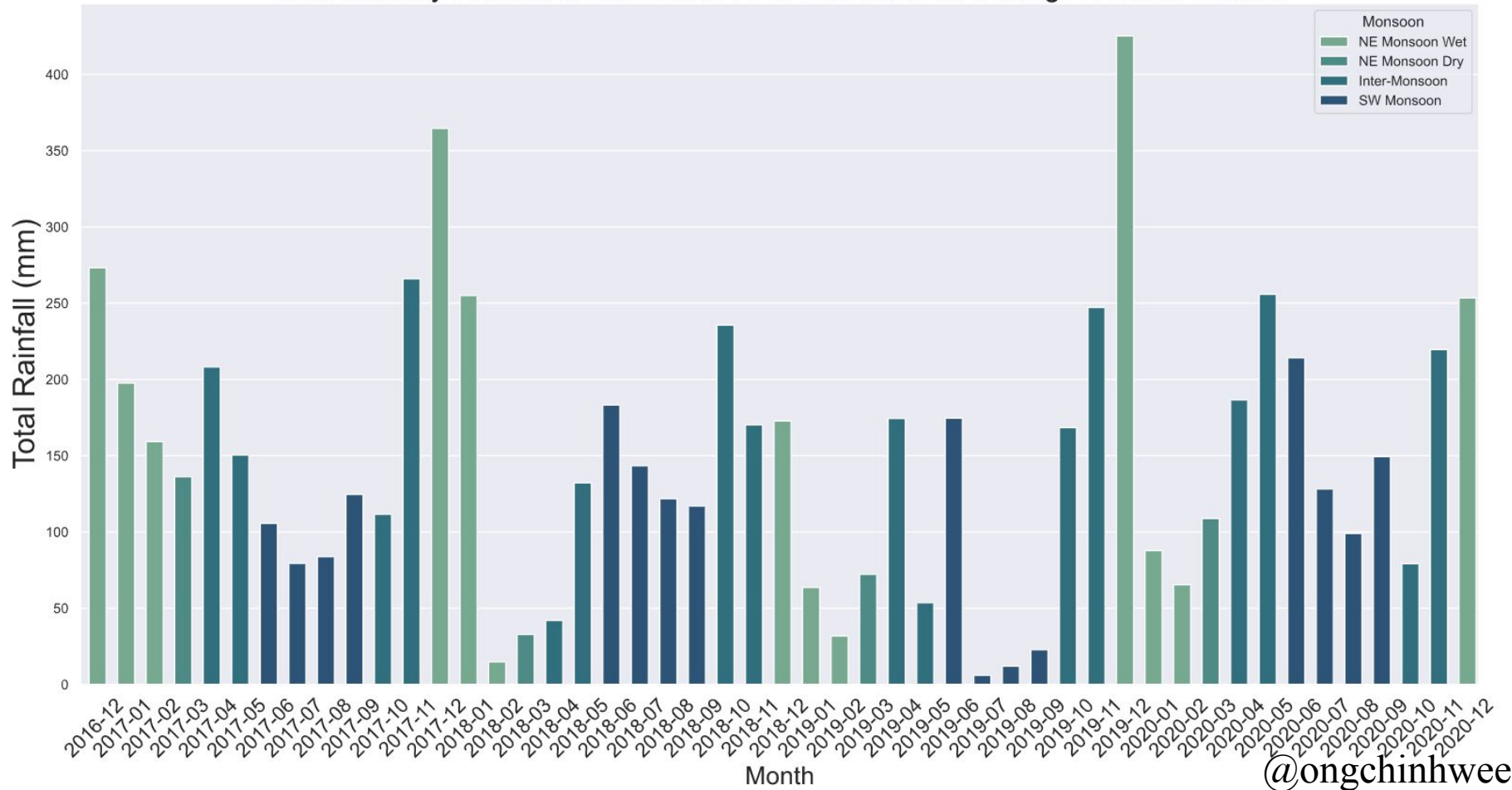
@ongchinwee

Total weekly sampling of rainfall at Changi Weather Station

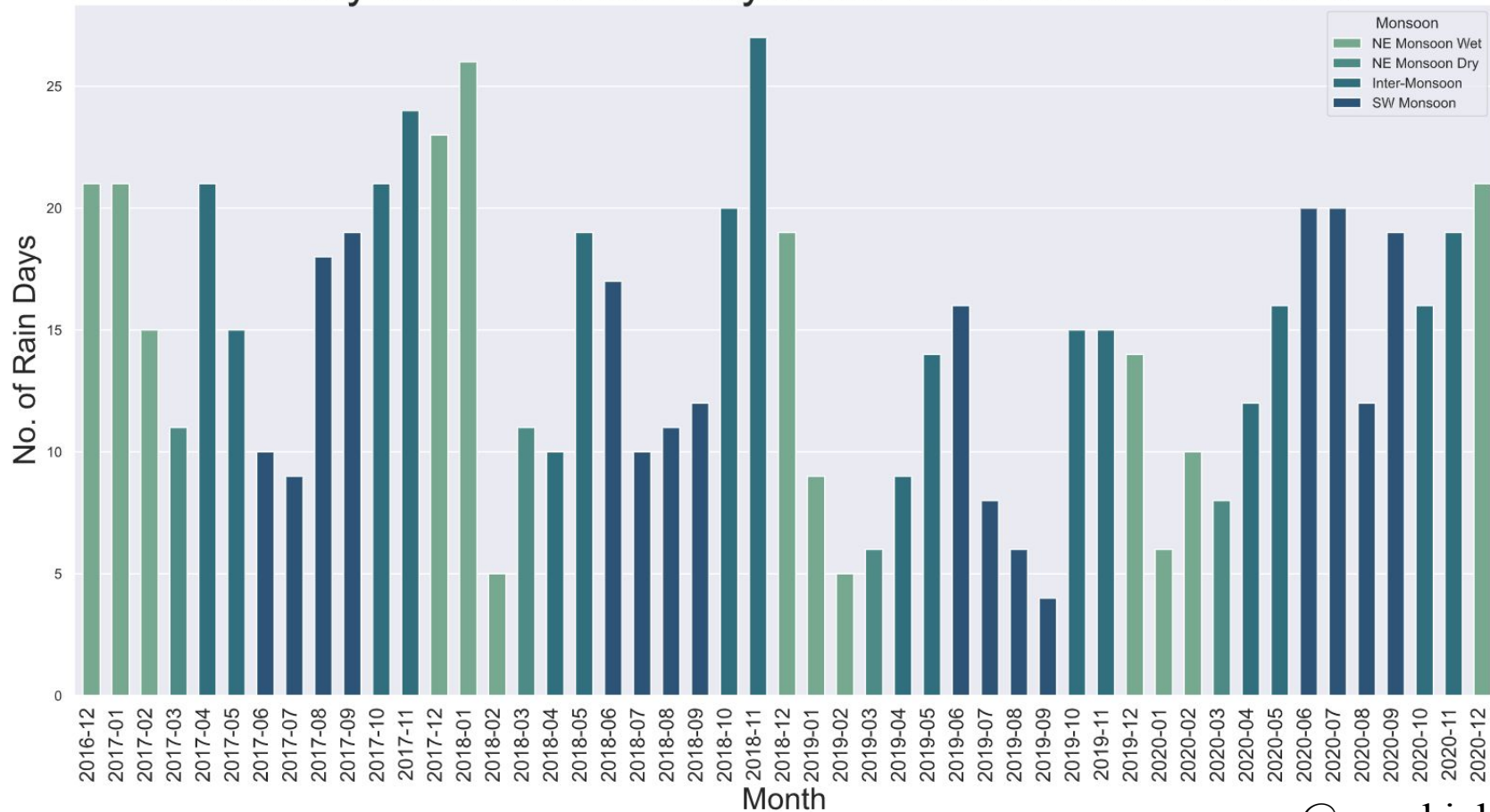




Total Monthly Rainfall from Dec 2016 to Dec 2020 at Changi Weather Station



# Monthly number of rain days from Dec 2016 to Dec 2020



# Time Series Analysis for Forecasting

Analyse and forecast time series using “statsmodels.tsa”

“statsmodels” library:

Python library for statistical models, tests and exploration

“statsmodels.tsa”:

Model classes and functions for Time Series Analysis

# Time Series Analysis for Forecasting

**Stationarity:** Stationary vs Non-Stationary

- Augmented Dickey-Fuller (ADF) Test

**Patterns:** Trend, Seasonality, Cycles (and Noise)

- Moving Averages
- STL Decomposition

**Autocorrelation:** Relationship between a time series and a lagged version of itself

# Augmented Dickey-Fuller (ADF) Stationary Test

```
from statsmodels.tsa.stattools import adfuller

def ADF_test(timeseries):
    dfctest = adfuller(timeseries.dropna(), autolag="AIC")
    print("Test statistic = {:.3f}".format(dfctest[0]))
    print("P-value = {:.3f}".format(dfctest[1]))
    print("Critical values :")
    for k, v in dfctest[4].items():
        print(
            f"\t{k}%: {v:.3f} - The data is {"not" if v < dfctest[0] else ""} stationary
            with {100 - int(k[:-1])}% confidence")
```

# Augmented Dickey-Fuller (ADF) Stationary Test

## Total Daily Rainfall

Test statistic = -5.710

P-value = 0.000

Critical values :

1%: -3.585 - The data is

**stationary with 99% confidence**

5%: -2.928 - The data is

stationary with 95% confidence

10%: -2.602 - The data is

stationary with 90% confidence

## Monthly Daily Rainfall

Test statistic = -13.590

P-value = 0.000

Critical values :

1%: -3.435 - The data is

**stationary with 99% confidence**

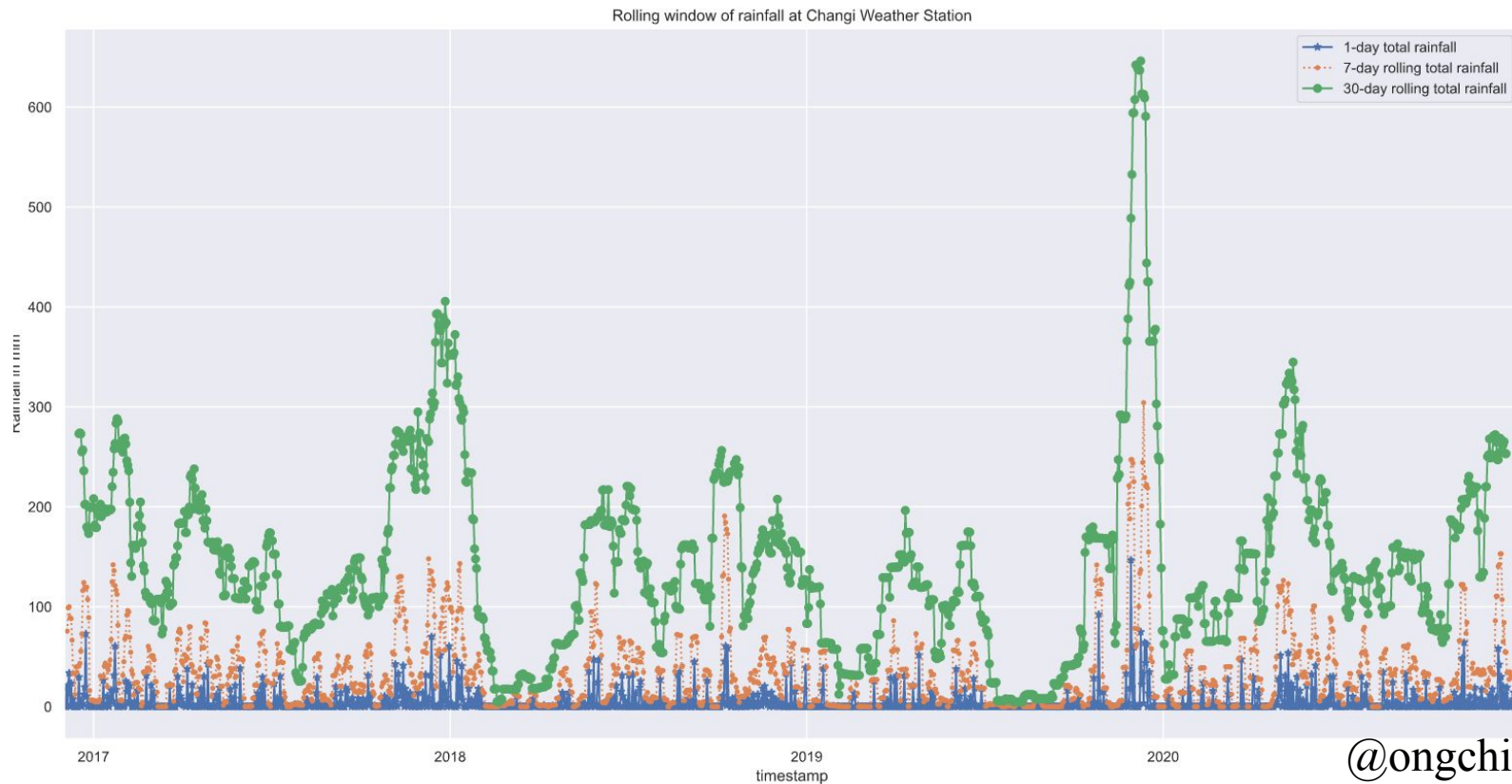
5%: -2.864 - The data is

stationary with 95% confidence

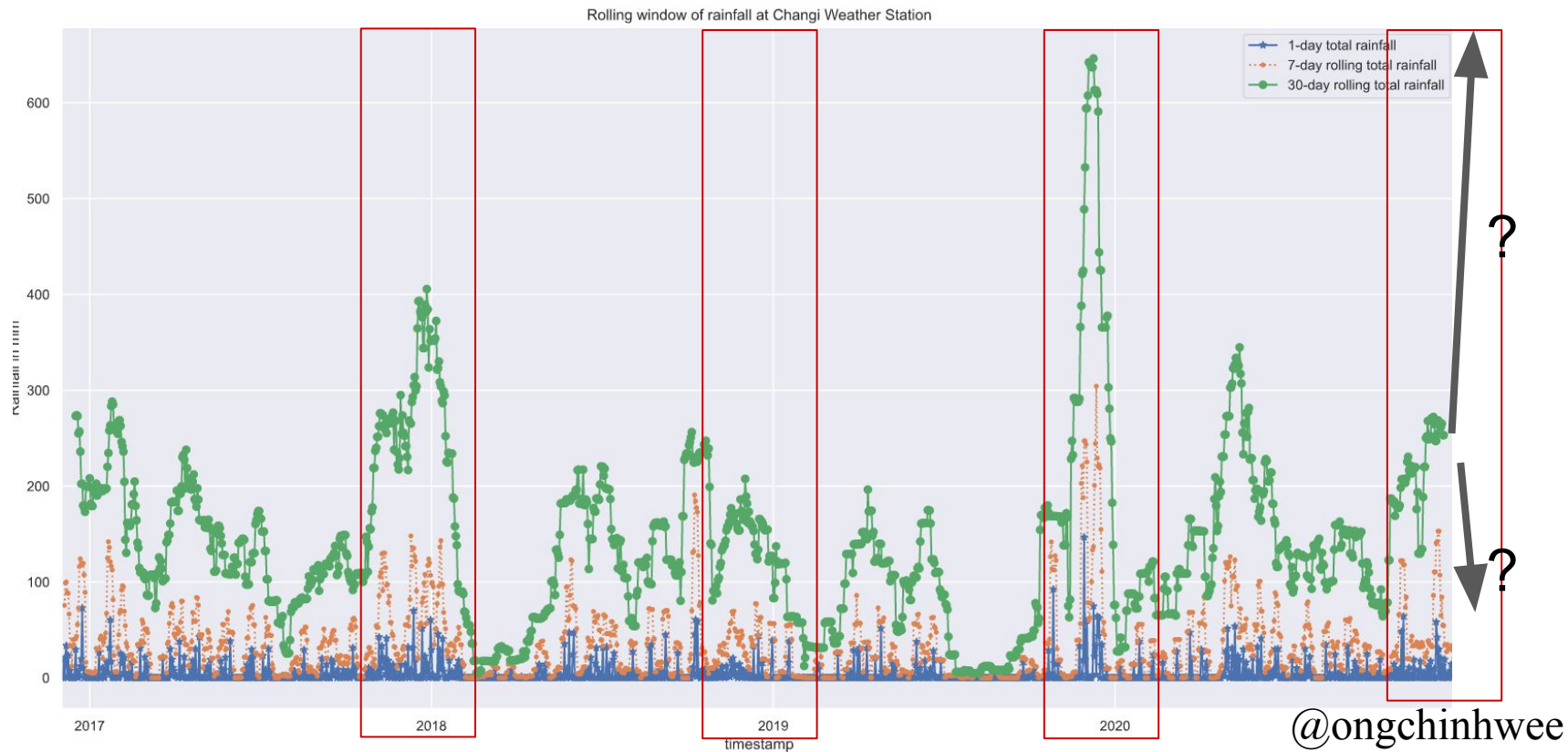
10%: -2.256 - The data is

stationary with 90% confidence

# Analyzing Rainfall with Moving Averages

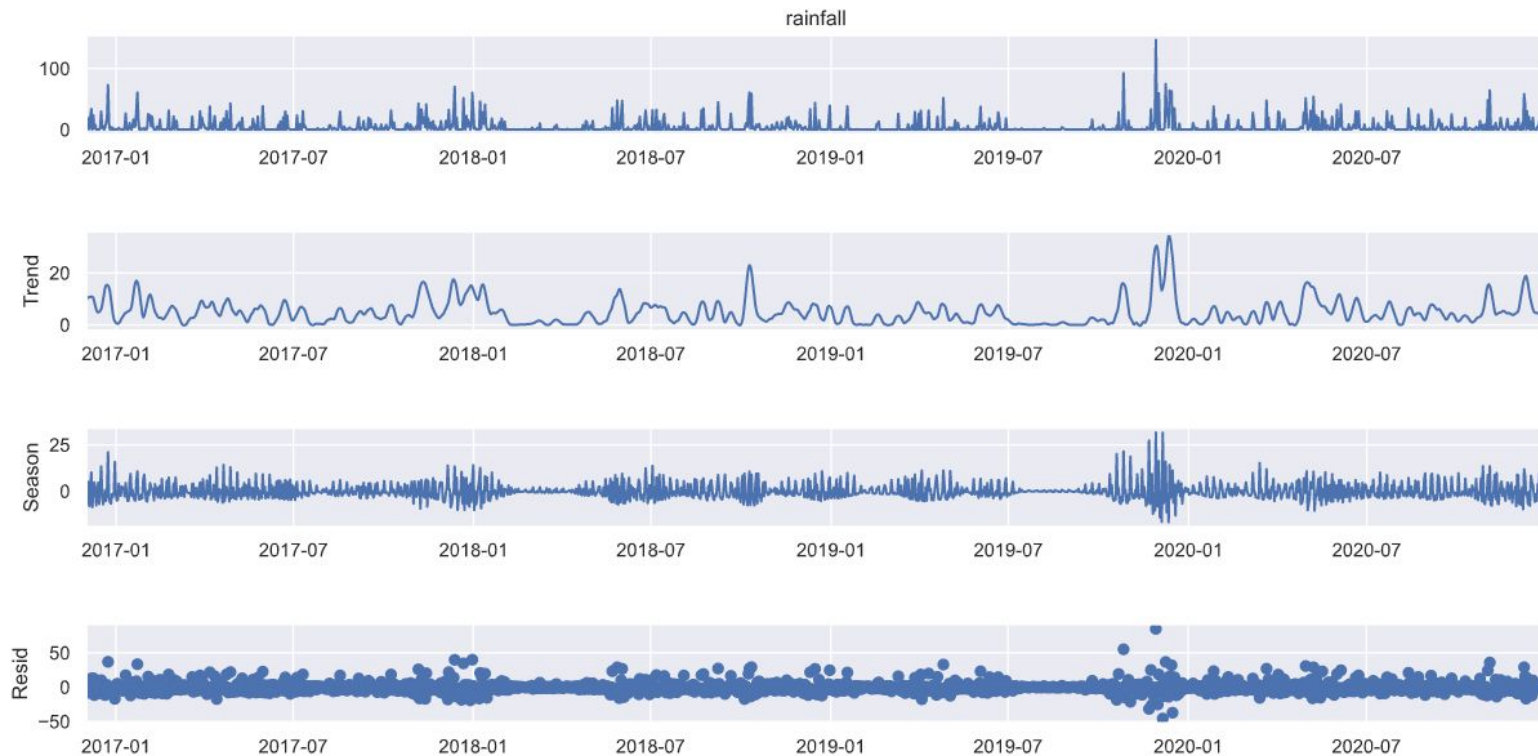


# Analyzing Rainfall with Moving Averages

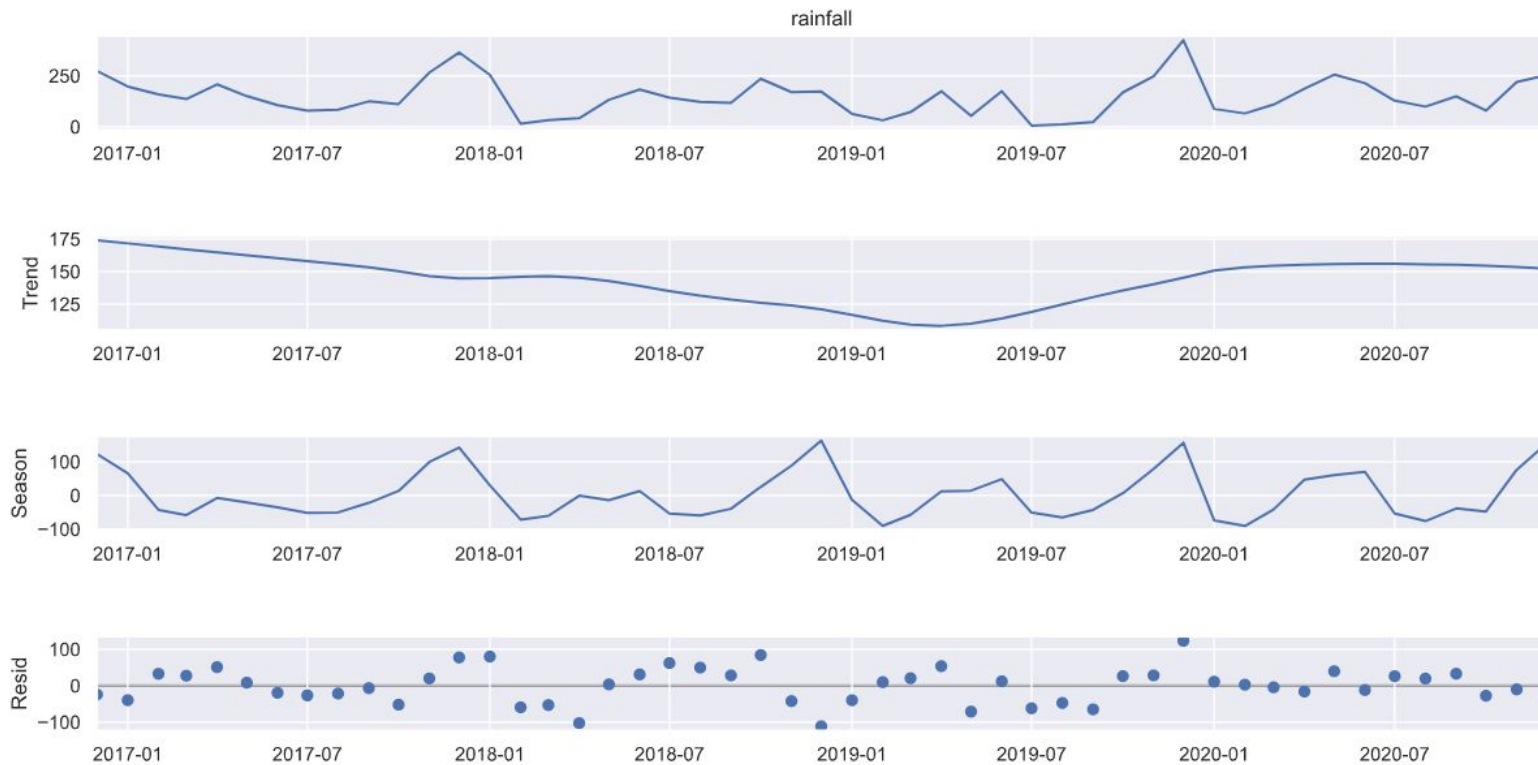




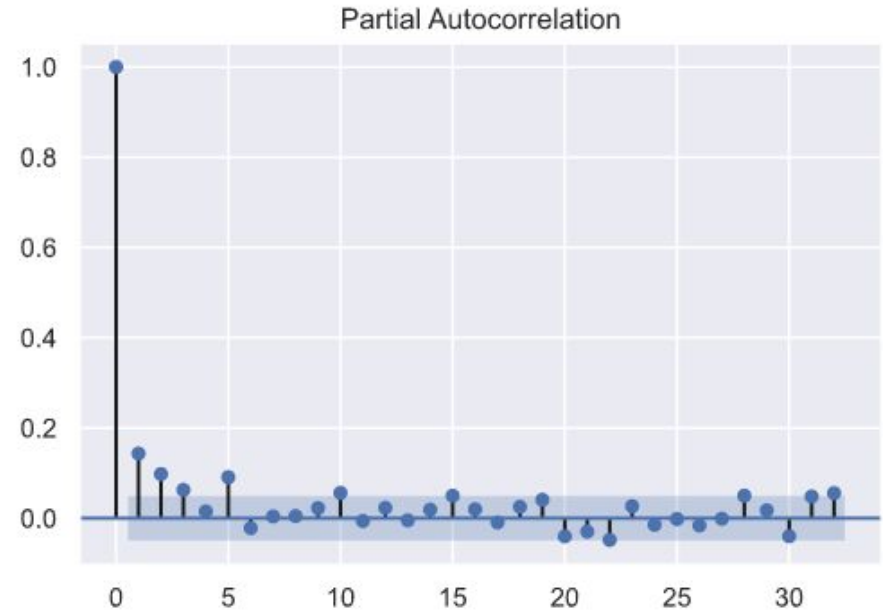
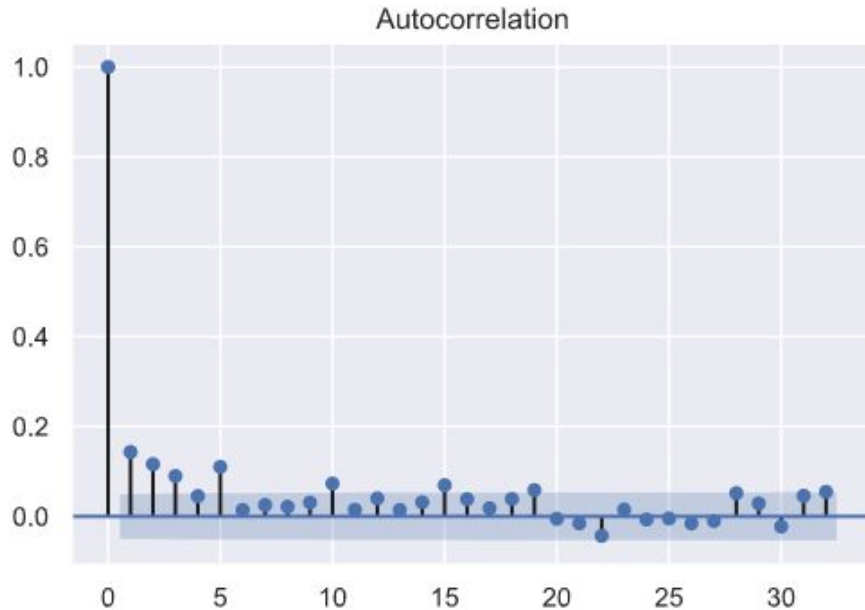
# STL Decomposition of Daily Rainfall



# STL Decomposition of Monthly Rainfall

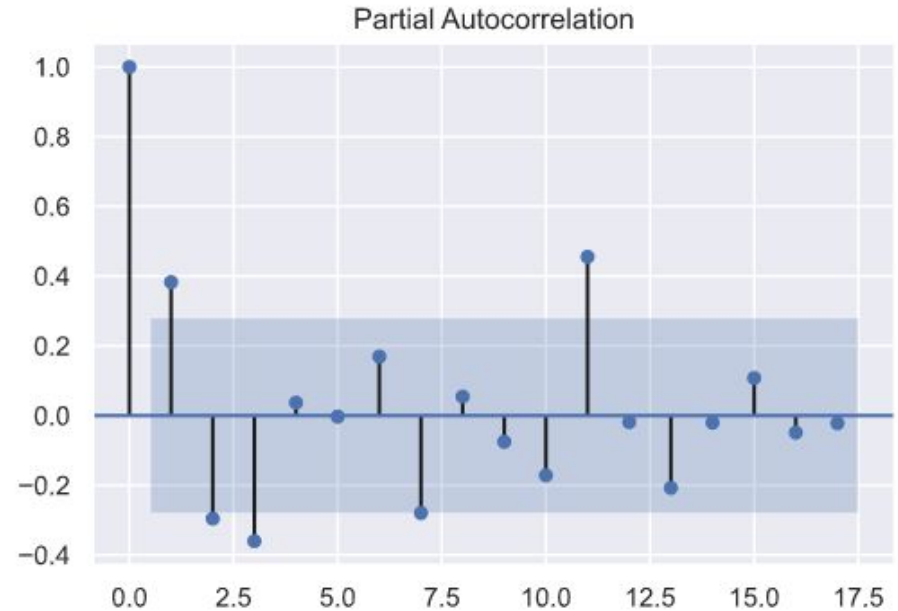
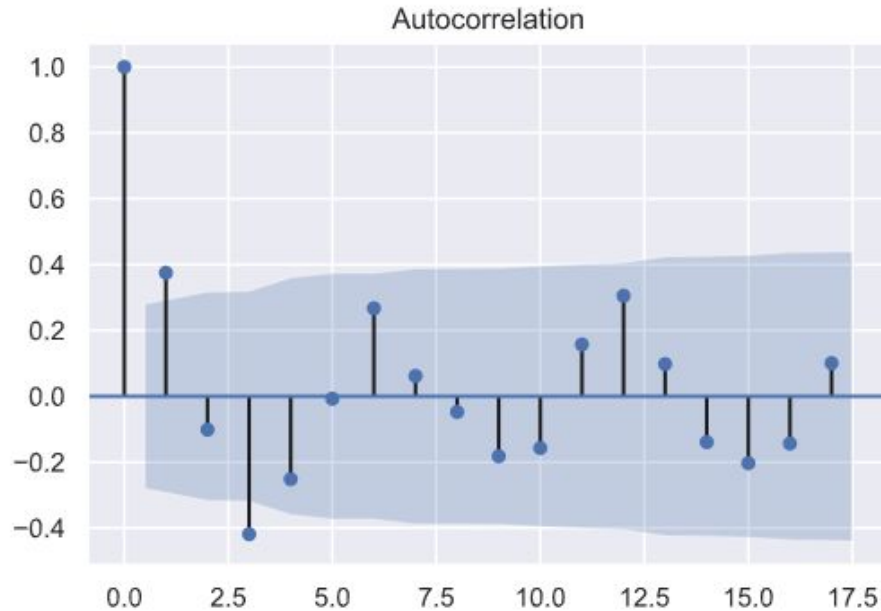


# Autocorrelation of Daily Rainfall



**Low correlation** between daily rainfall and its own lagged values.

# Autocorrelation of Monthly Rainfall



Most positive: **1st coefficient ( $r_1$ )**; Most negative: **3rd coefficient ( $r_3$ )**

**Recap:**

Could we **predict heavier rainfall**  
with weather data?

# Rainfall Forecasting with ARIMA models

## ARIMA(p,d,q) model

(AutoRegressive Integrated Moving Average)

where:

p: order of the **autoregressive** part;

d: degree of first **differencing** involved;

q: order of the **moving average** part.

# Rainfall Forecasting with ARIMA models

1. Apply **rolling forecast technique** with ARIMA( $p, d, q$ ) on time series data
2. Minimise **root-mean-squared-error (RMSE)**
3. Use optimized order parameters ( $p, d, q$ ) to run **rolling forecast for next  $N$  cycles**
  - a. Daily Forecast:  $N = 61$
  - b. Monthly Forecast:  $N = 13$

# Forecasting of Daily Rainfall with ARIMA

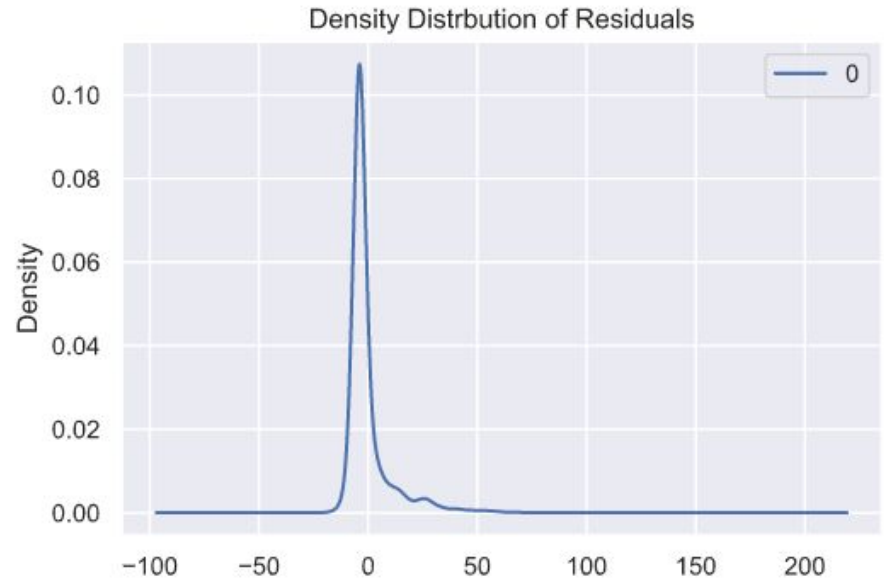
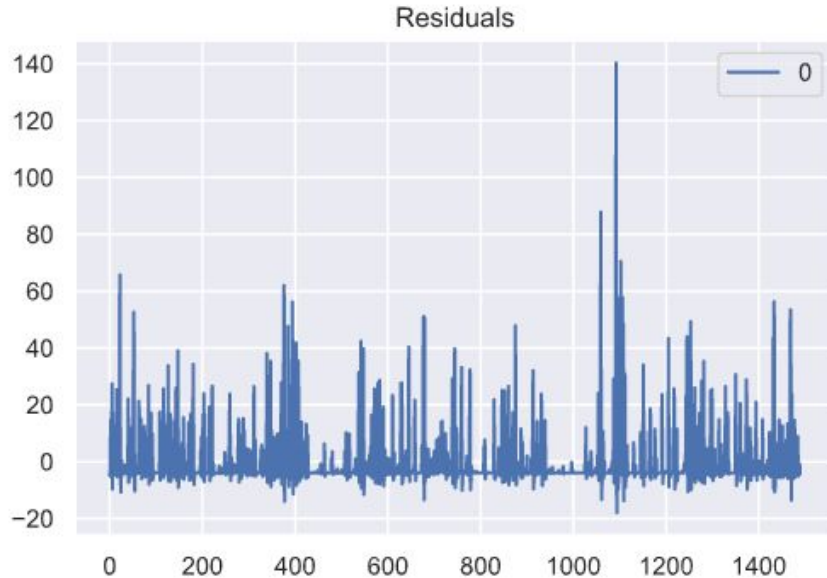
## SARIMAX Results

```
=====
Dep. Variable:          y      No. Observations:          1490
Model:                ARIMA(0, 0, 2)    Log Likelihood      -5687.965
Date:                Thu, 14 Jan 2021    AIC                  11383.930
Time:                23:48:12           BIC                  11405.156
Sample:              0              HQIC                  11391.840
                        - 1490
Covariance Type:      opg
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
const          4.8190        0.542        8.890      0.000        3.757        5.881
ma.L1          0.1197        0.020        5.900      0.000        0.080        0.159
ma.L2          0.0972        0.018        5.291      0.000        0.061        0.133
sigma2        121.0694        1.816       66.679      0.000       117.511       124.628
=====
Ljung-Box (L1) (Q):          0.06    Jarque-Bera (JB):          51835.43
Prob(Q):                    0.81    Prob(JB):              0.00
Heteroskedasticity (H):      1.66    Skew:                  4.14
Prob(H) (two-sided):        0.00    Kurtosis:             30.68
=====
```



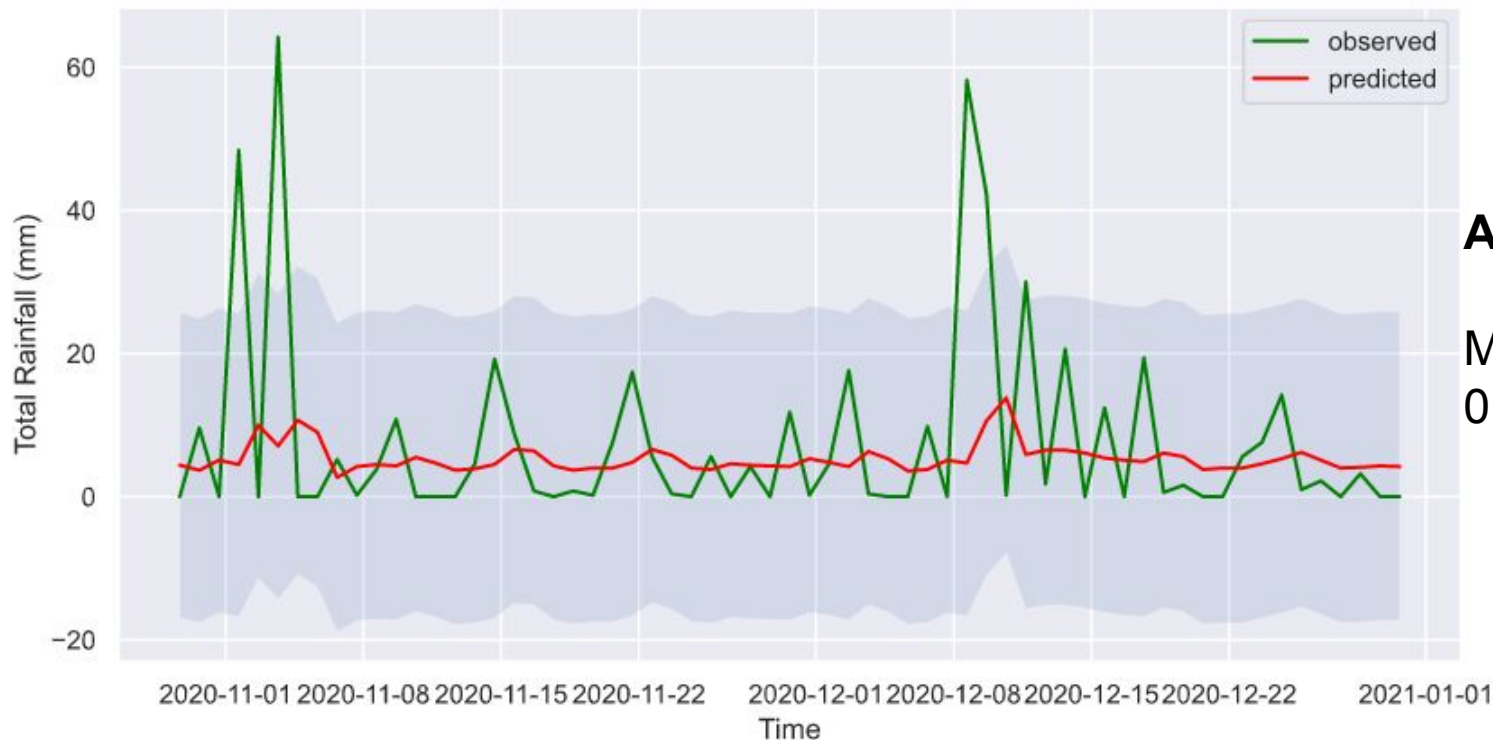
# Residual Errors for ARIMA Daily Rainfall Model

**ARIMA(0,0,2)**



# Forecasting of Daily Rainfall with ARIMA

Predicted vs Observed Rainfall



**ARIMA(0,0,2)**

Min-max error:  
0.757

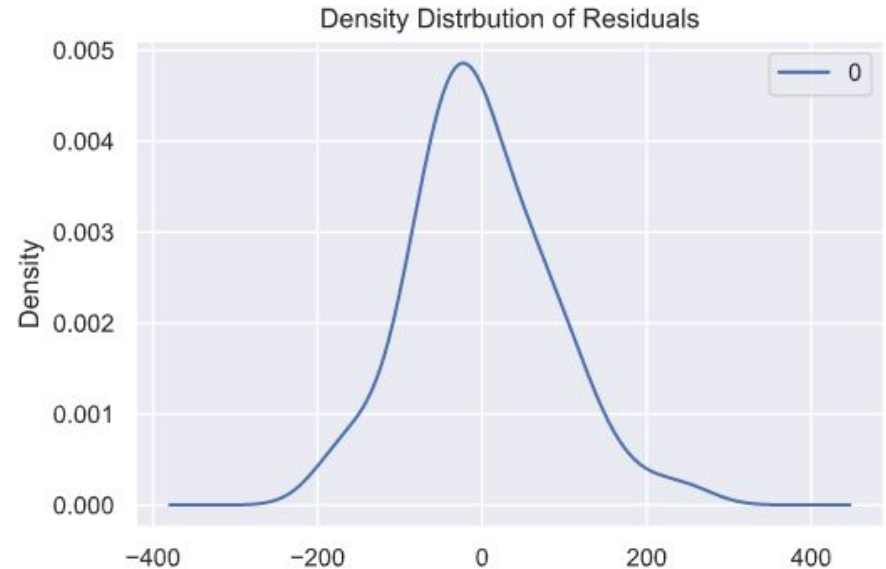
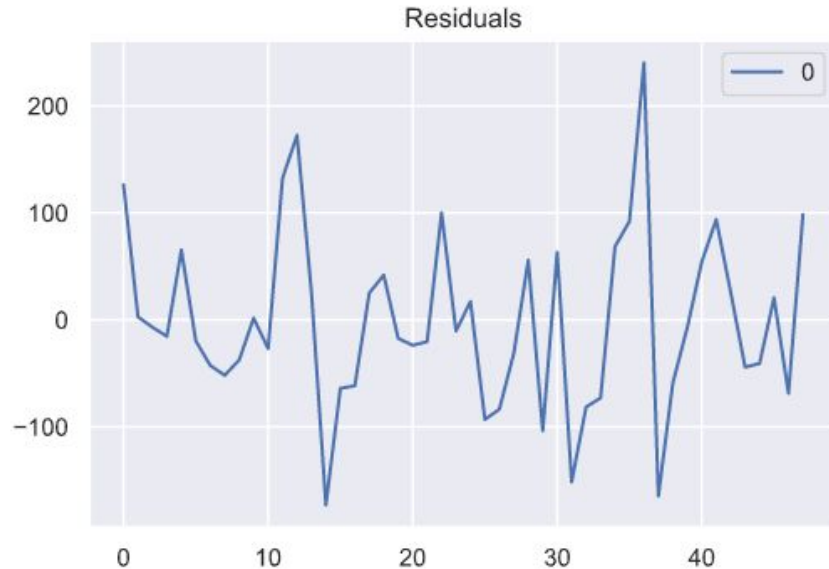
# Forecasting of Monthly Rainfall with ARIMA

## SARIMAX Results

```
=====
Dep. Variable:          y      No. Observations:          48
Model:                 ARIMA(1, 0, 0)      Log Likelihood      -280.016
Date:                 Thu, 14 Jan 2021      AIC                  566.033
Time:                 23:50:13      BIC                  571.646
Sample:                0      HQIC                  568.154
                        - 48
Covariance Type:        opg
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
const         146.8506      21.313      6.890      0.000      105.077      188.624
ar.L1           0.3792       0.116      3.269      0.001       0.152       0.607
sigma2        6810.5983    1297.728      5.248      0.000     4267.097     9354.099
=====
Ljung-Box (L1) (Q):          0.42      Jarque-Bera (JB):          1.77
Prob(Q):                     0.52      Prob(JB):              0.41
Heteroskedasticity (H):      1.35      Skew:                  0.41
Prob(H) (two-sided):         0.56      Kurtosis:              3.48
=====
```

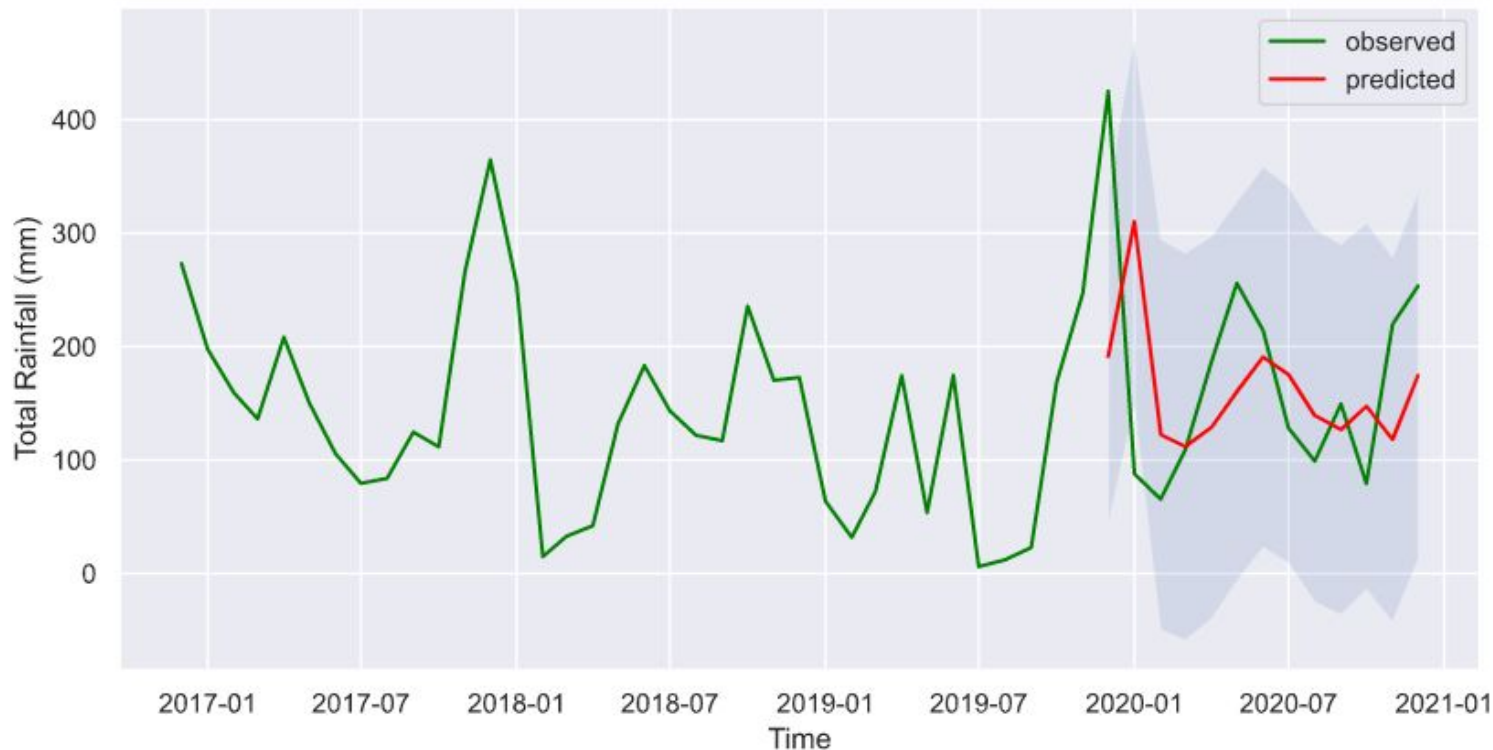
# Residual Errors for ARIMA Monthly Rainfall Model

**ARIMA(1,0,0)**



# Forecasting of Monthly Rainfall with ARIMA

Predicted vs Observed Rainfall



**ARIMA(1,0,0)**

Min-max error:  
0.346

@ongchinwee

# Key Takeaways

- With **climate change**, rainfall patterns are becoming **more extreme** and **more challenging to predict**
  - **Highest rainfall** in December 2019 (NE Monsoon)
  - **Higher-than-expected rainfall** in May 2020 (Inter-Monsoon) - also **earlier-than-expected monsoon**
- Rainfall data from weather station + ARIMA **may not be sufficient enough** to predict more “erratic” spikes in daily rainfall



# Reach out to me!



: ongchinhwee



: @ongchinhwee



: hweecat



: <https://ongchinhwee.me>

And check out my project  
on:



hweecat/api-scraping-nea-datasets