# Air Quality - Forecast

SARA SANCHEZ

TIME SERIES ANALYSIS AND MODELING

# Objective

To predict the Relative Humidity (RH) based on several features from a data set over the period March 2004 to February 2005.
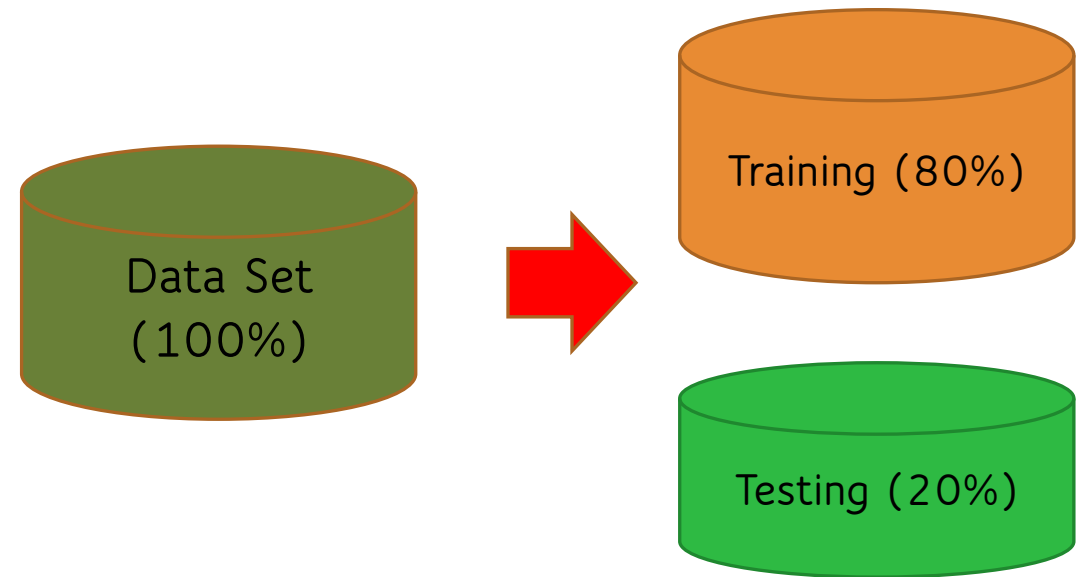
# Data Set

• Repository: UCI Machine Learning

• 9,358 Instances

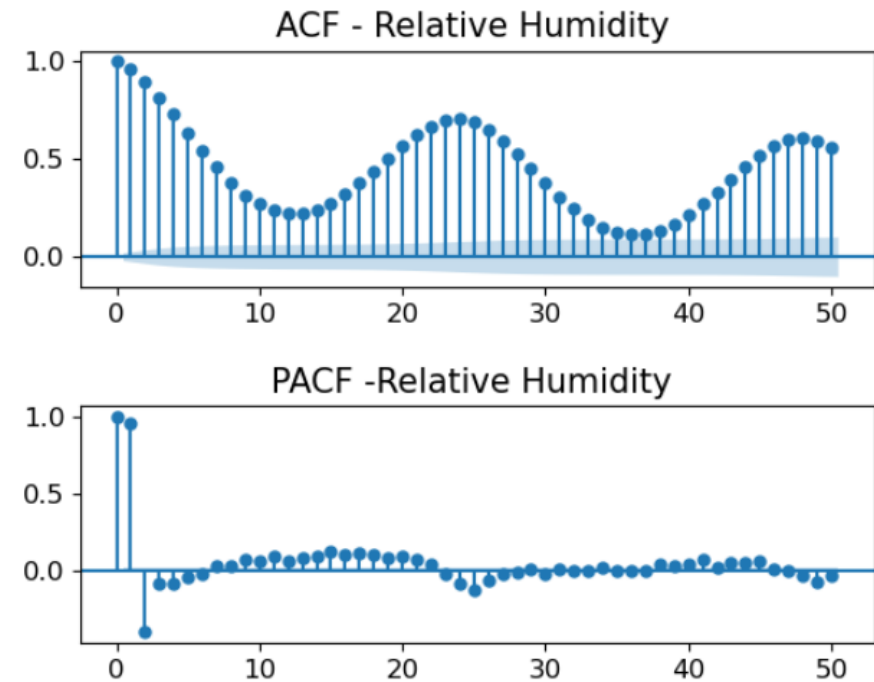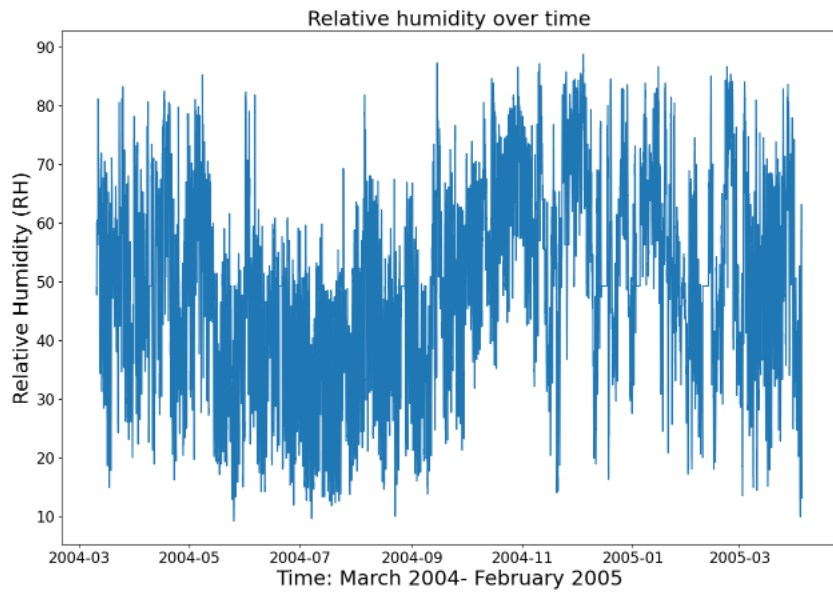• Dependent variable: Relative Humidity (RH)

# Data Preprocesing

As parte of the data preprocesing:

- Removing Unamed columns

- Manage the dates

- Changing some datatypes
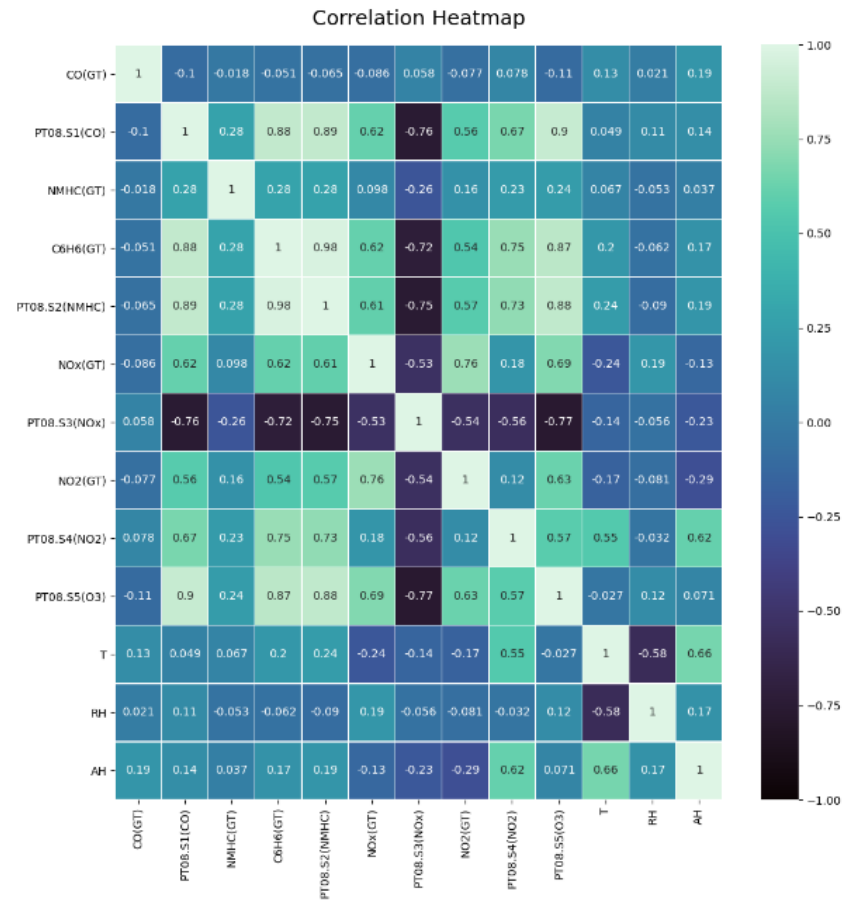
- Handling Null values

- Handling NaN values

Data Set
(100%)

Training (80%)

Testing (20%)

# EDA

# Correlation Heatmap

# Checking - Stationality

```
ADF Statistic: -7.391164
p-value: 0.000000
Critical Values:
    1%: -3.431
    5%: -2.862
    10%: -2.567
```

```
Results of KPSS Test:
Test Statistic            2.963095
p-value                   0.010000
LagsUsed                 52.000000
Critical Value (10%)      0.347000
Critical Value (5%)       0.463000
Critical Value (2.5%)     0.574000
Critical Value (1%)       0.739000
dtype: float64
```
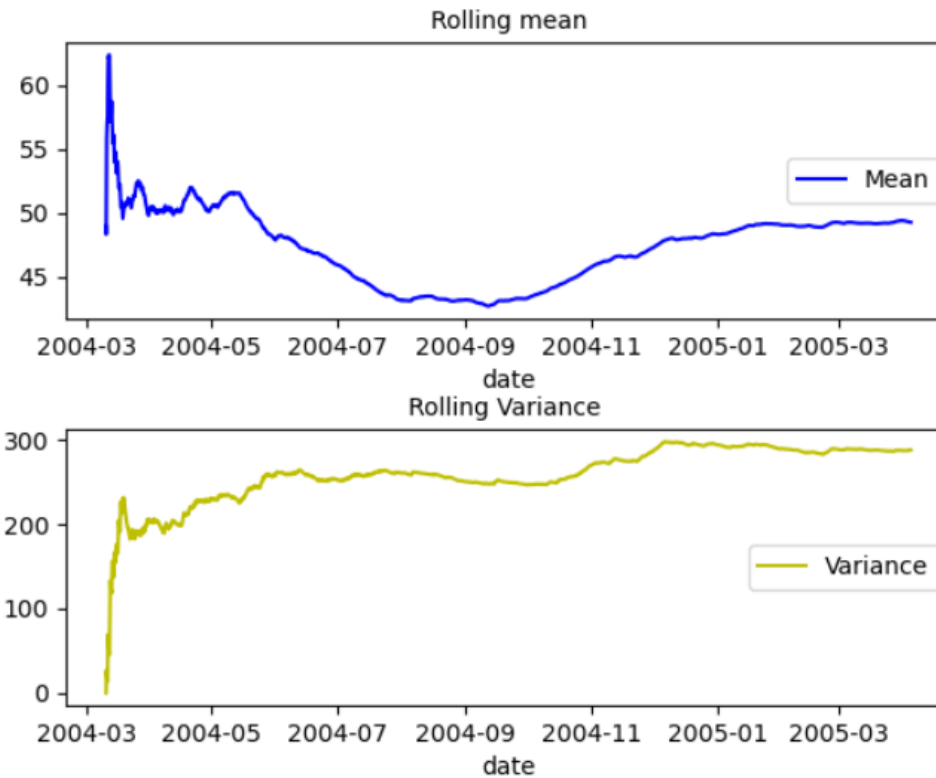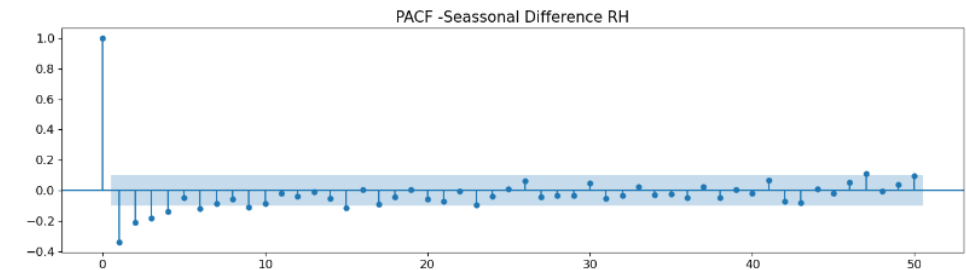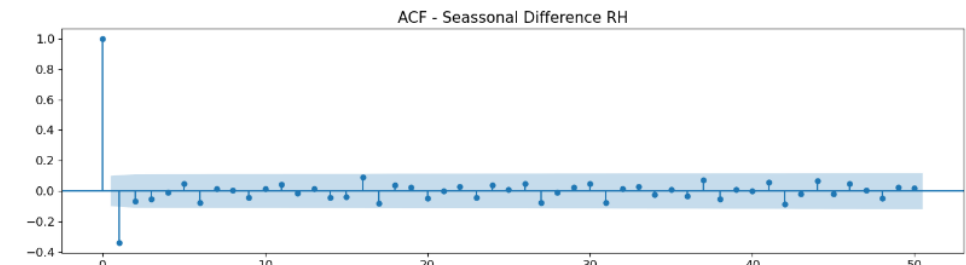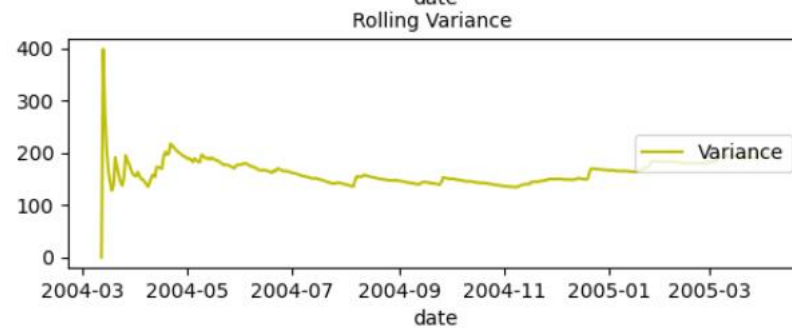


Rolling mean



Rolling Variance

# Seasonal difference – 24 periods

# Time Series Decomposition



Trend, Seasonality, Residual components using STL Decomposition



Original vs Seasonally adjusted

Strength of trend for Air quality dataset is 0.879

Strength of seasonality for Air quality dataset is 0.807

# Holt-Winter

# Feature Selection

# Multiple Linear Regression

# ARMA - Models



GPAC Table for RH

ARMA (1,0)

ARMA (2,1)

# ARMA (1,0)



ARMA Model Results

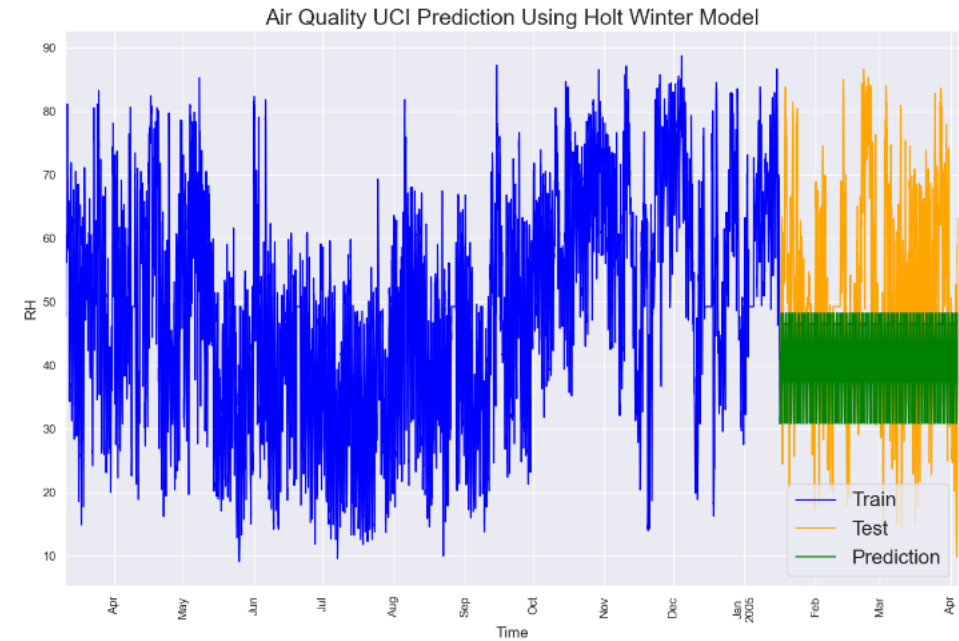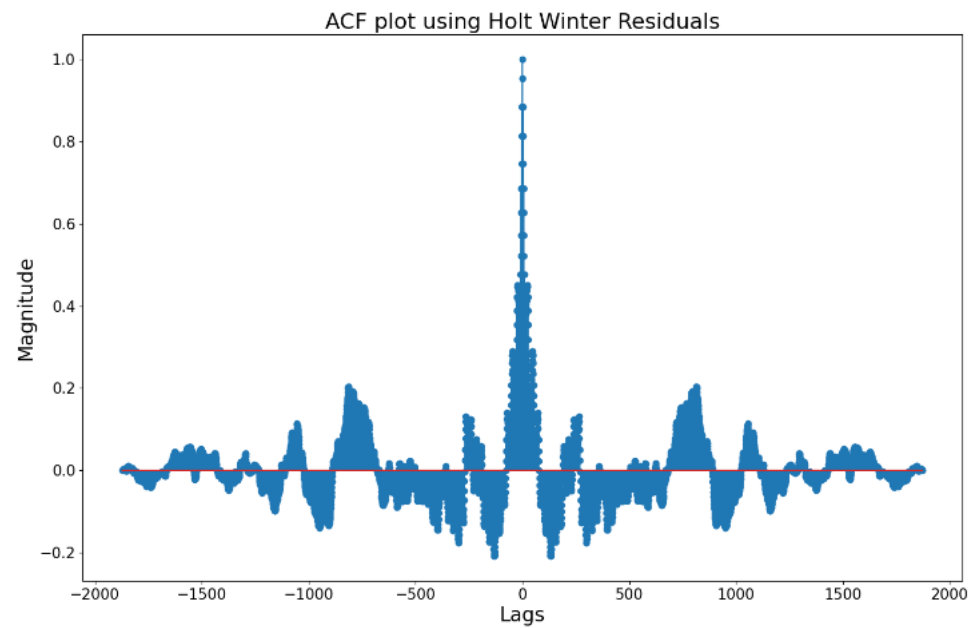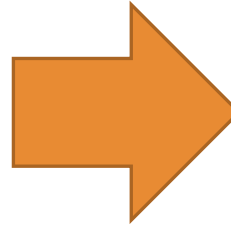| | | | |
|---|---|---|---|
| Dep. Variable: | RH | No. Observations: | 9357 |
| Model: | ARMA(1, 0) | Log Likelihood | -27529.364 |
| Method: | css-mle | S.D. of innovations | 4.586 |
| Date: | Wed, 04 May 2022 | AIC | 55062.729 |
| Time: | 00:04:44 | BIC | 55077.016 |
| Sample: | 03-10-2004 | HQIC | 55067.581 |
| | - 04-04-2005 | | |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| ar.L1.RH | 0.9629 | 0.003 | 345.193 | 0.000 | 0.957 | 0.968 |

Roots

| | Real | Imaginary | Modulus | Frequency |
|---|---|---|---|---|
| AR.1 | 1.0385 | +0.0000j | 1.0385 | 0.0000 |

## LM – Parameter Estimation

PARAMETER ESTIMATED
===============================================================:
LM - The AR coefficient a0 is: 0.9960090800559419
The AR coefficient a0 is: 0.995995649538989

# ARMA (2,1)

```
                    ARMA Model Results
==============================================================================
Dep. Variable:                  RH   No. Observations:                 9357
Model:                  ARMA(2, 1)   Log Likelihood              -26603.418
Method:                    css-mle   S.D. of innovations              4.154
Date:            Wed, 04 May 2022    AIC                          53214.837
Time:                    00:12:15    BIC                          53243.412
Sample:                  03-10-2004  HQIC                         53224.542
                       - 04-04-2005
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1.RH       1.6266      0.019     85.550      0.000       1.589       1.664
ar.L2.RH      -0.6685      0.018    -36.681      0.000      -0.704      -0.633
ma.L1.RH      -0.3351      0.025    -13.540      0.000      -0.384      -0.287
                                Roots
==============================================================================
                  Real          Imaginary           Modulus         Frequency
------------------------------------------------------------------------------
AR.1            1.2166           -0.1258j            1.2230           -0.0164
AR.2            1.2166           +0.1258j            1.2230            0.0164
MA.1            2.9844           +0.0000j            2.9844            0.0000
------------------------------------------------------------------------------
```



ACF plot for ARMA(2,1) Residuals

## LM – Parameter Estimation

LM - The AR coefficient a0 is: 1.5249896053962666
LM - The AR coefficient a1 is: -0.5298550819910275
LM - The MA coefficient b0 is: -0.1843455766129109
The AR coefficient a0 is: 1.5190248961591442
The AR coefficient a1 is: -0.5239233881532351
The MA coefficient b0 is: -0.17136892003716545



Air Quality UCI Prediction Using ARMA(2, 1) Model

# SARIMA (0,0,0) x(0,1,1,24)



```
                              SARIMAX Results
==========================================================================
Dep. Variable:                    RH   No. Observations:           7485
Model:            SARIMAX(0, 1, [1], 24)   Log Likelihood       -28612.806
Date:                  Wed, 04 May 2022   AIC                   57229.612
Time:                          00:42:17   BIC                   57243.447
Sample:                       03-10-2004   HQIC                  57234.365
                            - 01-16-2005
Covariance Type:                    opg
==========================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
--------------------------------------------------------------------------
ma.S.L24      -0.6081      0.008    -73.520      0.000      -0.624      -0.592
sigma2       125.2868      1.544     81.151      0.000     122.261     128.313
==========================================================================
Ljung-Box (L1) (Q):              6459.24   Jarque-Bera (JB):        892.57
Prob(Q):                            0.00   Prob(JB):                  0.00
Heteroskedasticity (H):             1.12   Skew:                      0.33
Prob(H) (two-sided):                0.00   Kurtosis:                  4.56
==========================================================================
```

# Model's Comparison

```
                                BASE MODEL COMPARISON
==================================================================================================================
                              Model         MSE        RMSE  Residual Mean  Residual Variance  Train Residual Mean  Train Residual Variance       Q Value
                     Average Model  261.847603   16.181706       1.231013        260.332210        -2.269931e-09              294.757911  310975.648032
                       Naive Model  592.261784   24.336428      18.218932        260.332210         1.698792e+01              294.757911  310975.648032
                       Drift Model  675.720932   25.994633      20.333690        262.261970         2.544018e+01              367.999923  305120.133891
  Simple Exponential Smoothing Model  582.260222   24.130869      17.942352        260.332210        -3.701496e-03               34.940240  310975.648032
                  Holt Winter Model  298.016539   17.263156       9.938193        199.248865        -5.502696e-02               19.478492  174779.944494
      Multiple Linear Regression Model   60.433018    7.773868      -0.640819         60.022369         1.103621e-01               33.677514   45559.267884
                    ARMA(1, 0) Model  986.902672   31.415007       2.426130        981.016567         2.185490e-03             1115.137897  316119.924675
                    ARMA(2, 1) Model 1003.068885   31.671263       2.389274        997.360256         9.273044e-04             1125.326443  315440.233711
  SARIMA (0, 0, 0) (0, 1, 1, 24) Model  418.858124   20.466024     -14.524584        207.894582         2.365511e-01              137.638637  178853.964200
==================================================================================================================
```
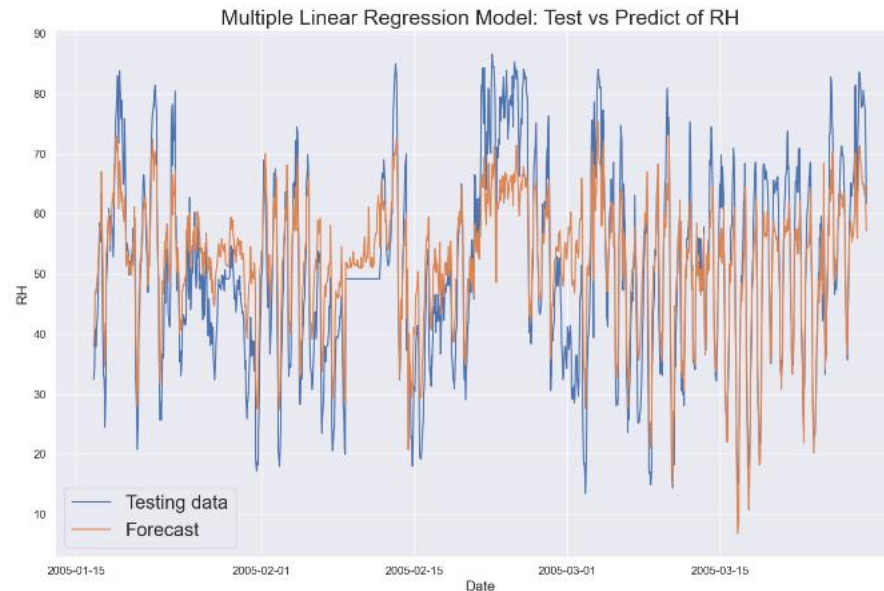
# Final Model – Multiple Linear Regression



Multiple Linear Regression Model: Test vs Predict of RH

Final Equation:

$$Y = 0.0006 * CD(GT) + 0.0135 * PT08.S1(CO) + 0.0033 * NMHC(GT) - 1.6980 * C6H6(G) + 0.0339 * PT08.S2(NMHC) + 0.0157 * NOx(GT) + 0.0104 * PT08.S3(N0x) - 0.0241 * NO2(GT) + 0.0122 * PT08.S4(NO2) - 0.0017 * PT08.S4(NO2) - 2.3287 * T + 34.5833 * AH$$

Thanks!!!!!!