

# Assignment-2-CH5440

CH18B020

March 13, 2022

## Question-1)

### Part a)

We have 4 independent variables: concentrations of  $CO_2$ ,  $CH_4$ ,  $N_2O$ ,  $O_3$ . And, we have one independent variable:  $T_{avg, deviation}$ . We can fit a model of the form  $X\beta + \beta_0 = y$  where we can estimate the parameters as,

- $\hat{\beta} = (X^T X)^{-1} X^T y$
- $\hat{\beta}_0 = \bar{y} - \bar{x}^T \hat{\beta}$

The obtained model is given as:

$$T_{deviation} = 11.8 + 0.0607x_1 + 0.00591x_2 - 0.14652x_3 + 0.00804x_4 \quad (1)$$

where,

1.  $x_1$  is concentration of  $CO_2$
2.  $x_2$  is concentration of  $CH_4$
3.  $x_3$  is concentration of  $N_2O$
4.  $x_4$  is concentration of  $O_3$

The temperature deviation is positively correlated with the concentration of all gases other than  $N_2O$ , for which it is negatively correlated. This is unexpected because the correlation estimate between that and temperature deviation turns out to be positive (approx. 0.88). The coefficient values were also verified using `fitlm()` function in MATLAB.

### Part b)

Confidence intervals:

| Term      | Lower Bound | Estimate | Upper Bound |
|-----------|-------------|----------|-------------|
| Intercept | -0.2598     | 11.7998  | 23.8594     |
| $x_1$     | 0.0336      | 0.0607   | 0.0878      |
| $x_2$     | 0.0039      | 0.0059   | 0.0079      |
| $x_3$     | -0.2208     | -0.1465  | -0.0722     |
| $x_4$     | -0.0025     | 0.0080   | 0.0186      |

If we keep the bound for the residuals as 2, we don't see any residual greater than that value as seen in figure 1. However, the residual corresponding to the 13th data point is somewhat off positioned.

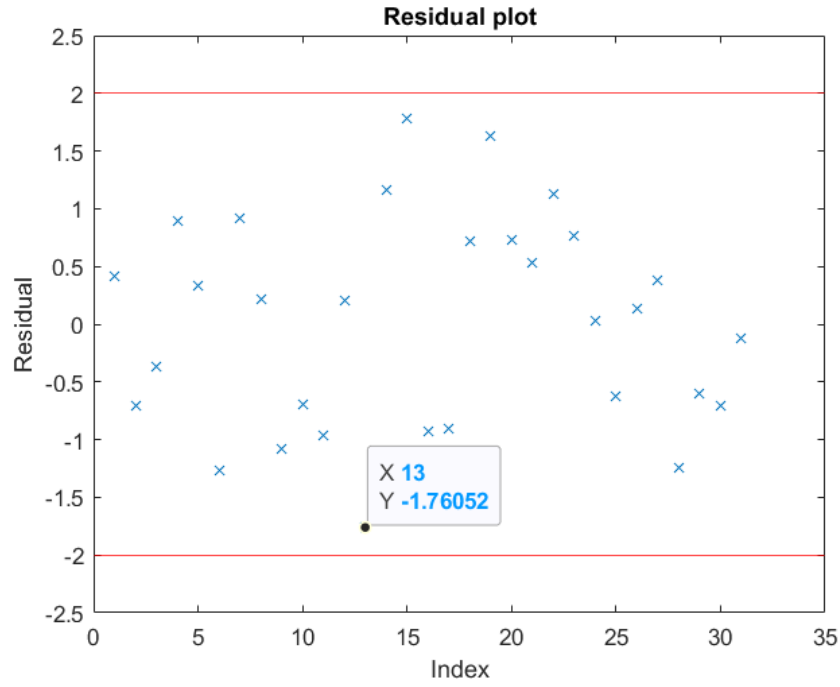


Figure 1: Residual plot for full data

Removing that point and retraining, the residuals are better now.

$$\hat{\beta}_{new} = [0.0631, 0.0066, -0.1573, 0.0086]^T \quad (2)$$

$$\hat{\beta}_{0,new} = 11.8075 \quad (3)$$

No outliers seen in the residuals of the new model.

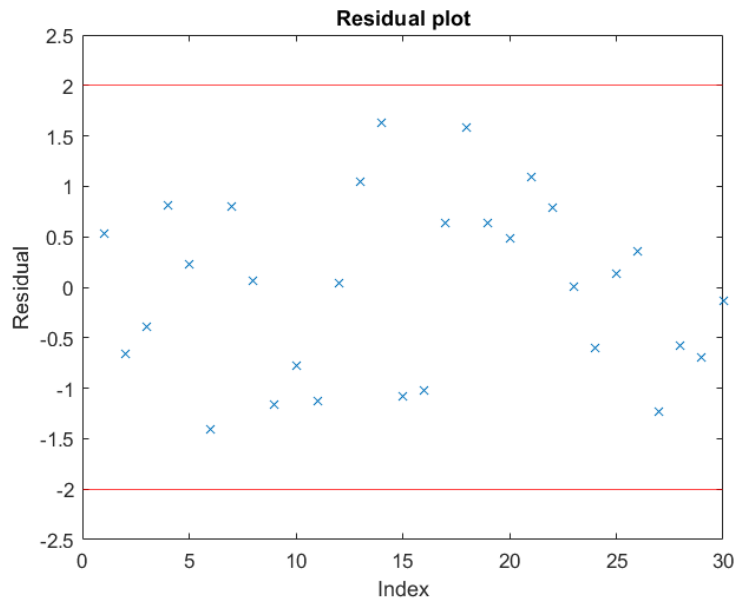


Figure 2: Residual plot after outlier removal

## Part c)

fitlm() function is used to find out the pValues and standard errors.

mdl\_c =

Linear regression model:

y ~ 1 + x1 + x2 + x3 + x4

Estimated Coefficients:

|             | Estimate  | SE        | tStat   | pValue    |
|-------------|-----------|-----------|---------|-----------|
|             | -----     | -----     | -----   | -----     |
| (Intercept) | 12.95     | 11.46     | 1.13    | 0.26919   |
| x1          | 0.063059  | 0.025766  | 2.4473  | 0.02176   |
| x2          | 0.0066088 | 0.0019334 | 3.4183  | 0.0021657 |
| x3          | -0.15733  | 0.07074   | -2.2241 | 0.035407  |
| x4          | 0.0086162 | 0.010032  | 0.85883 | 0.39859   |

Number of observations: 30, Error degrees of freedom: 25

Root Mean Squared Error: 0.122

R-squared: 0.858, Adjusted R-Squared: 0.835

F-statistic vs. constant model: 37.7, p-value = 3.02e-10

We can see that Pvalue for concentration of Ozone coefficient is high ( $> 0.05$ ) meaning it is insignificant. So we drop that first and rebuild a model.

mdl\_c1 =

Linear regression model:

y ~ 1 + x1 + x2 + x3

Estimated Coefficients:

|             | Estimate  | SE        | tStat   | pValue    |
|-------------|-----------|-----------|---------|-----------|
|             | -----     | -----     | -----   | -----     |
| (Intercept) | 17.266    | 10.247    | 1.6849  | 0.10397   |
| x1          | 0.067786  | 0.025044  | 2.7067  | 0.011845  |
| x2          | 0.0064572 | 0.0019155 | 3.3709  | 0.0023508 |
| x3          | -0.16813  | 0.069261  | -2.4274 | 0.022435  |

Number of observations: 30, Error degrees of freedom: 26

Root Mean Squared Error: 0.122

R-squared: 0.854, Adjusted R-Squared: 0.837

F-statistic vs. constant model: 50.5, p-value = 5.51e-11

We can see that Pvalue for intercept is high ( $> 0.05$ ) meaning it is insignificant. So we drop that and rebuild the model.

mdl\_c2 =

Linear regression model:

y ~ x1 + x2 + x3

Estimated Coefficients:

|    | Estimate  | SE        | tStat   | pValue     |
|----|-----------|-----------|---------|------------|
|    | -----     | -----     | -----   | -----      |
| x1 | 0.025869  | 0.0029784 | 8.6857  | 2.6637e-09 |
| x2 | 0.0043522 | 0.0015007 | 2.9001  | 0.0073299  |
| x3 | -0.052118 | 0.0077743 | -6.7039 | 3.3882e-07 |

Number of observations: 30, Error degrees of freedom: 27

Root Mean Squared Error: 0.126

We can see that Pvalue for all coefficients is very low ( $< 0.05$ ) which means all coefficients are highly significant.

### Part d)

GWP of the gases is simply ratio of its regression coefficient (adjusted for units) with the regression coefficient of  $CO_2$ .

1. GWP of  $CO_2 = 1$  (by definition)

2. GWP of  $CH_4 = \frac{0.0059}{0.0607} * 10^3 = 97.405$

3. GWP of  $N_2O = \frac{-0.1465}{0.0607} * 10^3 = -2413.3$

We observe that  $CH_4$  GWP is close to the values observed over a 20 year horizon, but  $N_2O$  is not.

## Question-2)

### Part a)

Considering,

$$y = \ln(P^{sat}) \quad (4)$$

$$x = \frac{1}{T} \quad (5)$$

we perform OLS and obtain,

$$A' = 4.7607 \quad (6)$$

$$B' = -37.896 \quad (7)$$

### Part b)

The optimization problem is set similar to OLS since measurements of  $y$  is noise-free, but of course, it is nonlinear in this case.

$$\min_{\hat{P}_1^{sat}, \dots, \hat{P}_{100}^{sat}, A, B, C} \sum_{i=1}^{100} (P_i^{sat} - \hat{P}_i^{sat})^2 \quad (8)$$

$$\text{s.t. } \ln(\hat{P}_i^{sat}) = A - \frac{B}{T_i + C} \quad (9)$$

$$(10)$$

We can eliminate the equality constraint by substituting  $P_i^{sat}$  back in the objective. After doing so, using `lsqnonlin()` to solve the problem, one obtains,

$$A = 14.1018 \quad (11)$$

$$B = 2821.4489 \quad (12)$$

$$C = 228.7554 \quad (13)$$

### Part c)

We are given that,

$$\sigma_{\epsilon_x} = 0.18 \quad (14)$$

$$\sigma_{\epsilon_y} = 2 \quad (15)$$

So we can use these values to set up a WTLS style optimization problem as given below.

$$\min_{\hat{P}_1^{sat}, \dots, \hat{P}_{100}^{sat}, \hat{T}_1, \dots, \hat{T}_{100}, A, B, C} \sum_{i=1}^{100} \frac{(P_i^{sat} - \hat{P}_i^{sat})^2}{\sigma_{\epsilon_y}^2} + \frac{(T_i - \hat{T}_i)^2}{\sigma_{\epsilon_x}^2} \quad (16)$$

$$\text{s.t. } \ln(P_i^{sat}) = A - \frac{B}{T_i + C} \quad (17)$$

$$(18)$$

To improve convergence, we have the initial guess for A,B,C as the solution of part b) (OLS problem), and the initial guess for temperature as the temperature measurement. We obtain the following estimates,

$$A = 14.1217 \quad (19)$$

$$B = 2835.2165 \quad (20)$$

$$C = 229.4130 \quad (21)$$

### Part d)

Listed below are the maximum absolute error in each case.

1. Part a): 39.8234
2. Part b): 4.4844
3. Part c): 4.2112

(3) a)  $S \mathbf{v}_1 = \lambda_1 \mathbf{v}_1$

$$\begin{bmatrix} 7 & 21 & 34 \\ 21 & 64 & 102 \\ 34 & 102 & 186 \end{bmatrix} \begin{bmatrix} v_{11} \\ v_{12} \\ v_{13} \end{bmatrix} = \begin{bmatrix} \lambda_1 v_{11} \\ \lambda_1 v_{12} \\ \lambda_1 v_{13} \end{bmatrix}$$

$$\Rightarrow (7 - \lambda_1) v_{11} + 21 v_{12} + 34 v_{13} = 0 \quad \text{--- (1)}$$

$$21 v_{11} + (64 - \lambda_1) v_{12} + 102 v_{13} = 0 \quad \text{--- (2)}$$

$$34 v_{11} + 102 v_{12} + (186 - \lambda_1) v_{13} = 0 \quad \text{--- (3)}$$

③  $\times$  ④ - ②

$$-34 v_{11} + (\lambda_1 + 1) v_{12} = 0$$

$$\Rightarrow v_{12} = \left( \frac{34}{\lambda_1 + 1} \right) v_{11} \quad \text{--- (4)}$$

$$- \frac{34 \times 21 \times 4}{21 \times 4} v_{12} = 3.012 v_{11}$$

34  $\times$  ② - 21  $\times$  ③

$$(34(64 - \lambda_1) - (21)(102)) v_{12}$$

$$+ (102 \times 34 - 21(186 - \lambda_1)) v_{13} = 0$$

$$\Rightarrow v_{13} = \frac{34\lambda_1 - 34}{21\lambda_1 - 438} v_{12} \quad \text{--- (5)}$$

$$\text{Product of roots} = |S|$$

$$= 7 \times 164 \times 186 - (102)^2$$

$$- 21(21 \times 186 - 34 \times 102)$$

$$+ 34(102 \times 21 - 64 \times 34)$$

$$= 146$$

$$\Rightarrow \lambda_2 \lambda_3 = \frac{146}{250 - 4} = 0.5831$$

Solving the quadratic

$$\lambda_2, \lambda_3 =$$

$$\frac{6.6 \pm \sqrt{(6.6)^2 - 4 \times 0.5831}}{2}$$

$$= 6.510, 0.0896$$

$$\therefore \lambda_2 = 6.510, \lambda_3 = 0.0896$$

Using (4) & (5),

$$V_{22} = \left( \frac{3\lambda_2}{\lambda_2 - 1} \right) V_{21}$$

$$3.544$$

$$\Rightarrow V_{22} = 2.16 V_{21}$$

$$V_{23} = \left( \frac{34 \times 11 - 234}{21 \lambda_1 - 438} \right) V_{22}$$

$$- 2.204$$

$$V_{23} = -0.6169 V_{21}$$

Normalizing,

$$V_2 =$$

$$\begin{bmatrix} 0.2330 \\ 0.8258 \\ -0.5135 \end{bmatrix}$$

$$\Rightarrow -v_{12} = \frac{-4820.4 - 8479.6}{-4820.4} v_{12}$$

$$\Rightarrow v_{13} = \frac{8479.6}{4820.4} \times \frac{250.4}{289.4} v_{11}$$

$$= 5.256 v_{11}$$

Normalising,

$$v_{11} = \frac{1}{\sqrt{1 + (3.012)^2 + (5.256)^2}} = 0.1629$$

$$v_{12} = 0.4906$$

$$v_{13} = 0.8560$$

$$\therefore v_1 = \begin{bmatrix} 0.1629 \\ 0.4906 \\ 0.8560 \end{bmatrix}$$

$$\text{Sum of eigen value} = \text{trace}(S) \\ = 7 + 64 + 186$$

$$\Rightarrow \lambda_1 + \lambda_2 + \lambda_3 = 257 \\ \Rightarrow \lambda_2 + \lambda_3 = 250 - 7 = 243$$



One again using ④ & ⑤,

$$V_{32} = \left( \frac{3 \times 0.0896}{0.0896 - 1} \right) \quad V_{31} = -0.2952 V_{31}$$

$$V_{33} = \left( \frac{34 \times 0.0896 - 34}{21 \times 0.0896 - 438} \right) V_{32}$$

$$= -0.0209 V_{32}$$

Normalising,

$$\therefore V_3 = \begin{bmatrix} 0.9589 \\ -0.2831 \\ -0.0200 \end{bmatrix}$$

Summary :

$$\lambda_1 = 250.4, \quad \lambda_2 = 6.5101, \quad \lambda_3 = 0.0896$$

$$V_1 = \begin{bmatrix} 0.1629 \\ 0.4906 \\ 0.8560 \end{bmatrix}, \quad V_2 = \begin{bmatrix} 0.2330 \\ 0.8258 \\ -0.5135 \end{bmatrix}, \quad V_3 = \begin{bmatrix} 0.9589 \\ -0.2831 \\ -0.0200 \end{bmatrix}$$

b) If we take just first component,

$$\% \text{ variance retained} = \frac{250.4}{250.4 + 6.5 + 0.0896} \times 100\%$$

$$= 97.436\% > 95\%$$

$\therefore$  It is enough if we retain just one principal component

c) If there are 2 linear relationships, then we just pick the eigen vectors corresponding to the 2 smallest eigenvalues

such that  $\underline{v} + \underline{z}_s = 0$   $\underline{z}_s$  is the shifted sample

$\underline{v}$  is the eigen vector.

So you get:

$$v_1(z_1 - \bar{z}_1) + v_2(z_2 - \bar{z}_2) + v_3(z_3 - \bar{z}_3) = 0$$

$$\bar{z}_1 = 9, \bar{z}_2 = 68, \bar{z}_3 = 129.$$

$$\Rightarrow v_1 z_1 + v_2 z_2 + v_3 z_3 - (v_1 \bar{z}_1 + v_2 \bar{z}_2 + v_3 \bar{z}_3) = 0$$

Substituting each of the eigen vectors,

Or.  $0.233 Z_1 + 0.8258 Z_2 - 0.5135 Z_3$   
 $+ 7.99 = 0$

$$\Rightarrow 0.233 Z_1 + 0.8258 Z_2 - 0.5135 Z_3 + 7.9901 = 0 \quad (6)$$

$$0.9589 Z_1 - 0.2831 Z_2 - 0.02 Z_3 + 13.201 = 0 \quad (7)$$

where  $Z_1$  is man

$Z_2$  is SVL

$Z_3$  is HLS

d) we need to project  $\begin{bmatrix} 10.1-9 \\ 73-68 \\ 135.5 \\ -129 \end{bmatrix}$  along the 3 different axes to get corresponding scores.

(mean shifted sample)

$$Z_5 = Z - \bar{Z} = \begin{bmatrix} 1.1 \\ 5 \\ 6.5 \end{bmatrix}$$

with  $Z_5$

Score along axis -1 =

$$v_1^T Z_5 = \boxed{8.18962}$$



$$\text{Score along axis-2} = V_2^T ZS = \boxed{1.04755}$$

$$\text{Score along axis-3} = V_3^T ZS = \boxed{-0.49071}$$

Since there are 2 linear relationships, we can assume the last 2 scores are just due to noise. The only score that matters for compression is along ~~axis~~ axis-1 with value which has a value of 8.1962

2) We have 2 linear relationships in (6) & (7) and we need to estimate 2 variables. So we can just solve the eqns to get  $\mu$  and  $\sigma$

Solving (6) & (7) for  $\mu$  by eliminating  $\sigma$ ,

$$\begin{aligned} & (6) \times (0.5135) + (7) \times (-0.0201) \\ \Rightarrow & (0.9589 \times 0.05135) Z_1 + (-0.283 \times 0.5135) Z_2 \\ & + (-0.0201)(0.233) Z_1 + (0.8259)(-0.0201) Z_2 = 0 \\ \Rightarrow & Z_1 Z_1 = \frac{-(0.8259 \times (-0.0201) - 4(0.233)(0.5135))}{0.9589 \times 0.05135 - 0.0201(0.233)} Z_2 \end{aligned}$$

$$\therefore \text{mass} = 19.2$$

$$0.4877 z_1 + 0.1618 z_2 + 6.6422 = 0$$

$$\Rightarrow \text{mass} = \frac{0.1618 \times 73 + 6.6184}{0.4877}$$

$$x_1 = 10.669$$

f) Here we have 2 equations and only one variable.  
So we can try to minimise the error as both equations

$$\min_{\hat{z}} (z - \hat{z})^+ (z - \hat{z})$$

Eliminating mass from ⑥ & ⑦,

$$0.8578 z_2 - 0.4877 z_3 + 4.588 = 0$$

$$2235$$

$$0.79185$$

We can set up a TLS problem so that

"  $z_2^+, z_3^+$  satisfies both ⑥ & ⑦

by obtaining  $\hat{z}_1^+$

$$\min_{z_2^*, z_3^*} (z_2 - z_2^*)^2 + (z_3 - z_3^*)^2$$

$$\text{s.t.} \begin{bmatrix} 0.8578 & -0.4877 \end{bmatrix} \begin{bmatrix} z_2^* \\ z_3^* \end{bmatrix} = -4.586$$

$$\Rightarrow z_2^* = \frac{-4.586 + 0.4877 z_3^*}{0.8578}$$

Subst. back we get an unconstrained quadratic problem.

~~$$\min_{z_2^*} (z_2 - z_2^*)^2 + (135.5 - z_3^*)^2$$~~

$$\min_{z_3^*} \left( 73 - \left( \frac{-4.586 + 0.4877 z_3^*}{0.8578} \right) \right)^2 + (135.5 - z_3^*)^2$$

$$\frac{\partial f}{\partial z_3} = 0 \Rightarrow 214.116 = 1.5685 z_3^*$$

$$\Rightarrow z_3^* = \frac{136.062}{1.5685}$$

$$z_2^* = \frac{72.0115}{1.5685}$$



$\therefore Z_1^+ \neq$  (from eqn (1) or (2))

$$= \boxed{10.34\% \text{ } g}$$