

Linköping University | Department of Computer and Information Science
Master's thesis, 30 ECTS | Statistics and Machine Learning
2021 | LIU-IDA/LITH-EX-A--2021/001--SE

Quantifying nitrogen oxides and ammonia via frequency modulation in gas sensors

- **DRAFT**

*Kvantifiering av kväveoxider och ammoniak via frekvensmodu-
lering i gassensorer*

Marcos Freitas Mourão dos Santos

Supervisor : Annika Tillander
Examiner : José M. Peña

External supervisor : Mike Andersson

Upphovsrätt

Detta dokument hålls tillgängligt på Internet - eller dess framtida ersättare - under 25 år från publiceringsdatum under förutsättning att inga extraordinära omständigheter uppstår.

Tillgång till dokumentet innebär tillstånd för var och en att läsa, ladda ner, skriva ut enstaka kopior för enskilt bruk och att använda det oförändrat för ickekommersiell forskning och för undervisning. Överföring av upphovsrätten vid en senare tidpunkt kan inte upphäva detta tillstånd. All annan användning av dokumentet kräver upphovsmannens medgivande. För att garantera äktheten, säkerheten och tillgängligheten finns lösningar av teknisk och administrativ art.

Upphovsmannens ideella rätt innehåller rätt att bli nämnd som upphovsman i den omfattning som god sed kräver vid användning av dokumentet på ovan beskrivna sätt samt skydd mot att dokumentet ändras eller presenteras i sådan form eller i sådant sammanhang som är kränkande för upphovsmannens litterära eller konstnärliga anseende eller egenart.

För ytterligare information om Linköping University Electronic Press se förlagets hemsida <http://www.ep.liu.se/>.

Copyright

The publishers will keep this document online on the Internet - or its possible replacement - for a period of 25 years starting from the date of publication barring exceptional circumstances.

The online availability of the document implies permanent permission for anyone to read, to download, or to print out single copies for his/hers own use and to use it unchanged for non-commercial research and educational purpose. Subsequent transfers of copyright cannot revoke this permission. All other uses of the document are conditional upon the consent of the copyright owner. The publisher has taken technical and administrative measures to assure authenticity, security and accessibility.

According to intellectual property law the author has the right to be mentioned when his/her work is accessed as described above and to be protected against infringement.

For additional information about the Linköping University Electronic Press and its procedures for publication and for assurance of document integrity, please refer to its www home page:
<http://www.ep.liu.se/>.

Abstract

The abstract resides in file `Abstract.tex`. Here you should write a short summary of your work.

Acknowledgments

Thank you for reading my draft! :)

Contents

Abstract	iii
Acknowledgments	iv
Contents	v
List of Figures	vii
List of Tables	viii
List of acronyms and abbreviations	ix
1 Introduction	1
1.1 Motivation	1
1.2 Aim	2
1.3 Research questions	3
2 Data	4
2.1 Data acquisition	4
2.2 Raw data	6
2.3 Pre-processing	7
3 Theory	10
3.1 Notation	10
3.2 Ordinary Least Squares Regression	11
3.3 Principal Component Analysis	11
3.4 Principal Component Regression	14
3.5 Partial Least Squares Regression	14
3.6 Ridge Regression	16
3.7 Cross Validation	17
4 Methods	19
4.1 Ordinary Least Squares	19
4.2 Principal Components Regression	19
4.3 Partial Least Squares Regression	20
4.4 Ridge Regression	20

4.5 Averaging features	20
5 Results	21
5.1 Individual exposures	21
5.2 Averaging features by gas mixture	26
6 Discussion	32
6.1 Results	32
6.2 Method	33
6.3 The work in a wider context	33
7 Conclusion	34
Bibliography	35
Appendix A Data acquisition time stamps	37

List of Figures

2.1	Schema of the data acquisition process.	4
2.2	An example of raw sensor response	5
2.3	Feature measurements times per cycle. The width of the red line indicates the duration of one of the feature measurement windows as an example.	6
2.4	A visualization of the feature measurement process.	6
2.5	Feature naming convention.	8
2.6	Pre-processed data structure.	8
4.1	A visualization of the feature averaging process.	20
5.1	Correlation matrix of features.	22
5.2	OLS prediction for individual exposures	22
5.3	PCA with two components for the three gases.	22
5.4	Explained variance	23
5.5	Number of PCs selection via CV.	23
5.6	PCR prediction for individual exposures	24
5.7	PLS with two components for the three gases.	24
5.8	PLS - Explained variance of X	25
5.9	Number of PLS components selection via CV.	25
5.10	PLSR prediction for individual exposures	26
5.11	CV of shrinkage factor λ	26
5.12	Shrinkage of ridge coefficients.	28
5.13	Ridge prediction for individual exposures	28
5.14	Slope features. Each line represents one of the 125 unique exposures.	29
5.15	Average features. Each line represents one of the 125 unique exposures.	29
5.16	OLS prediction for unique gas mixtures.	30
5.17	PCR prediction for unique gas mixtures.	30
5.18	PLSR prediction for unique gas mixtures.	31
5.19	Ridge prediction for unique gas mixtures.	31

List of Tables

2.1	Data acquisition details	5
2.2	Raw data column details	7
2.3	Sample of raw data.	7
2.4	Sample of pre-processed data.	9
5.1	Sample of averaged mixture data.	27
5.2	Sample of only average features.	27
A.1	Data acquisition timestamps.	37

List of acronyms and abbreviations

AC Alternating Current. 2

CV Cross Validation. vii, 17, 20, 23, 25, 26, 33

GBCO Gate Bias Cycled Operation. 2

Hz Hertz. 5, 7

mA milliamperes. 4

MSE Mean Squared Error. 17

NIPALS Nonlinear Iterative Partial Least Squares. 13, 15

OLS Ordinary Least Squares. v, 8, 11, 14, 15, 19, 21, 32, 33

PC Principal Component. vii, 12, 13, 14, 17, 20, 21, 23, 33

PCA Principal Components Analysis. vii, 11, 12, 13, 14, 19, 21, 22, 23, 33

PCR Principal Components Regression. v, 11, 14, 19, 23, 33

PLS Partial Least Squares. vii, 14, 15, 17, 20, 24, 25

PLSR Partial Least Squares Regression. v, vii, 2, 11, 14, 15, 20, 24, 26, 33, 34

ppm parts per million. 4, 7

RMSE Root Mean Squared Error. 17, 20, 23, 24, 33, 34

RSS Residual Sum of Squares. 11, 16

SAS Sensor and Actuator Systems. 4

SCR Selective Catalytic Reduction. 1, 2

SiC-FET Silicon Carbide Field Effect Transistor. 2, 4

TCO Temperature Cycled Operation. 2, 7, 33



1 Introduction

1.1 Motivation

Nitric Oxide (NO) and Nitrogen Dioxide (NO₂), commonly referred together as NO_x, are hazardous gases to the environment and to humans. Its main sources are combustion processes in transportation, and industrial processes such as (but not limited to) auto mobiles, trucks, boats, industrial boilers, turbines, etc. (USEPA 2019).

NO_x exposure to humans can cause respiratory illnesses such bronchitis, emphysema and can worsen heart disease (Boningari and Smirniotis 2016). Environmentally, NO_x are deemed precursors of adverse phenomena such as smog, acid rain, and the depletion of ozone (O₃) (Alberto Bernabeo, Webster, and Onofri n.d.). It is of high interest, therefore, to reduce NO_x emissions.

One well studied and successful method of reducing emissions is Selective Catalytic Reduction (SCR), which consists in the reduction of NO_x by ammonia (NH₃) into nitrogen gas (N₂) and water (H₂O) (Forzatti 2001), both harmless components. The process is based in the following reactions (Forzatti 2001):

- 4 NH₃ + 4 NO + O₂ → 4 N₂ + 6 H₂O
- 2 NH₃ + NO + NO₂ → 2 N₂ + 3 H₂O
- 8 NH₃ + 6 NO₂ → 7 N₂ + 12 H₂O

One key element in these reactions, however, is the amount of ammonia dosed into the SCR systems. Ammonia itself is hazardous to humans, causing skin and respiratory irritation, among other illnesses (ASTDR 2004). More importantly, ammonia is one of the main sources of nitrogen pollution and it has direct negative impact on biodiversity via nitrogen deposition in soil and water (Guthrie, Giles, Dunkerley, Tabaqchali, Harshfield, Ioppolo, and Manville 2018).

Hence it is also desired to keep ammonia emissions to a minimum. Too much ammonia in the SCR catalyst will guarantee NO_x reduction at the expense of undesired ammonia emissions. Concurrently, too little ammonia will impede SCR to occur properly, beating the purpose of the catalyst and as a consequence, undesired NO_x emissions.

To monitor gasses concentrations, chemical sensors are deployed, one of which is the Silicon Carbide Field Effect Transistor (SiC-FET). The identification and quantification of gasses is normally achieved through multiple sensor in so called sensor arrays. Ideally each sensor in the array needs to have different responses to different compounds (Bastuck 2019). The deployment of multiple sensors, on the other hand, proves itself cumbersome due to the increased chances of failure, and decalibration of the system should one or multiple sensors be replaced (Bastuck 2019).

One solution to this problem is the cycled operation of one single sensor, referred as virtual multi-sensor (Bastuck 2019). By cycling the working point parameters of the sensor, different substances react differently in the sensor surface, which in turn produces different responses. Temperature Cycled Operation (TCO), Gate Bias Cycled Operation (GBCO), and the combination of the two have been proven to increase selectivity of SiC-FET sensors (Bastuck 2019).

TCO, in contrast with a constant temperature evaluation, produces unique transient sensor responses, i.e. each gas mixture yields a slightly different sensor output. This unique gas signature increases selectivity (Bur, Bastuck, Lloyd Spetz, Andersson, and Schütze 2014). Additionally, the high temperatures reached in these cycles help in the cleansing of the sensor surface, preparing it for the new mixtures to come.

Frequency modulation tries to achieve the same goal: avoid steady state responses in exchange of unique signatures that could help identify/quantify the gasses at hand. It consists on operating the sensor in Alternating Current (AC). One then can regulate the frequency of this operation and create cycles of different frequencies, similar to what is done in TCO. This is equivalent to GBCO, but with more frequency changes and achieving overall higher frequencies.

The main question is: given these set of unique sensor responses, how one can quantify the gasses that produced them? Multivariate regression techniques have been shown to be successful: Partial Least Squares Regression (PLSR) has been used in chemometrics extensively and it has been proven to be good at this task (Bastuck 2019) (Wold, Sjöström, and Eriksson 2001). Other multivariate regression methods, naturally, can also be used.

1.2 Aim

The aim of this thesis is to investigate if frequency modulation can be used to simultaneously quantify the concentrations of NO_x and Ammonia in a particular gas mixture.

1.3 Research questions

1. Can frequency modulation be used to simultaneously quantify NO_x and Ammonia concentrations?
2. Which method yields best predictions of gas concentrations?

2 Data

2.1 Data acquisition

The data was acquired at the Sensor and Actuator Systems (SAS) laboratory at Linköping University. The experiment — as shown on Figure 2.1 — consisted of exposing different gas combinations to two SiC-FET sensors under a particular frequency cycle and recording its response, measured in milliamperes (mA). This is then used to extract secondary features, namely average and slope values from certain regions of the frequency cycle.

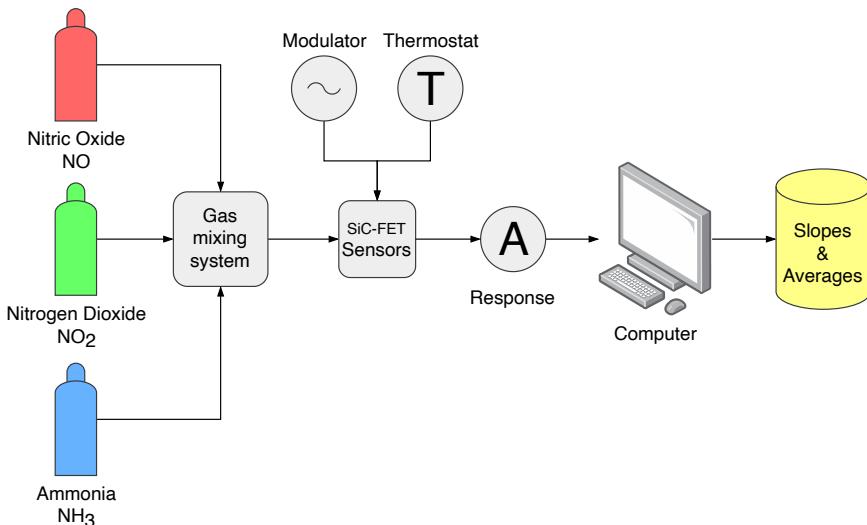


Figure 2.1: Schema of the data acquisition process.

In more detail, NO, NO₂ and NH₃ had five possible concentration values each: 5, 10, 20, 40, and 80 parts per million (ppm). The experiment was designed to encompass all possible combinations of these gases, amounting to 125 different gas mixtures. Each feature was sub-

mitted to the same frequency cycle four times. The cycle consists of 16 unique frequencies: 0.05, 0.1, 0.25, 0.5, 1, 2, 5, 10, 25, 50, 100, 200, 500, 1000, 2500 and 5000 Hertz (Hz). A typical raw sensor response for frequency modulation experiments is shown on Figure 2.2.

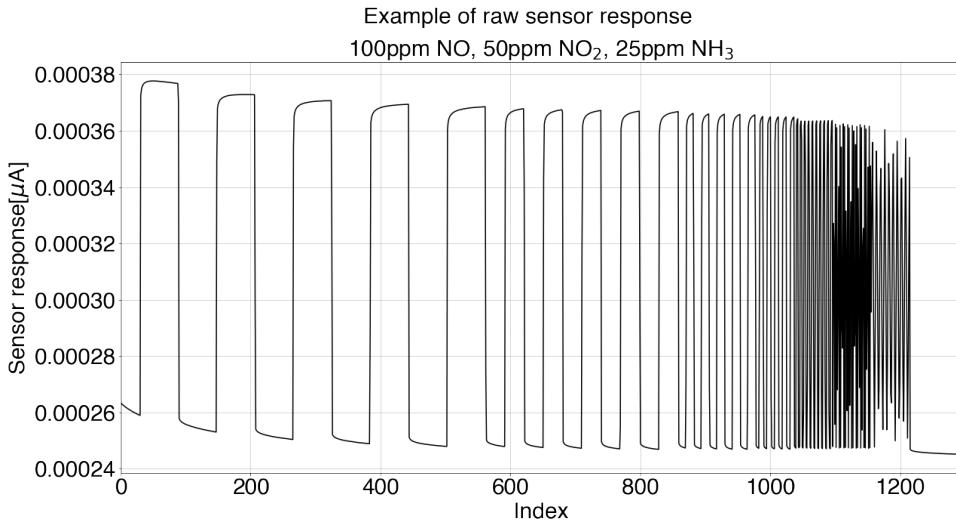


Figure 2.2: An example of raw sensor response

Throughout one cycle, several slope and average features were extracted. The sample rate for feature extraction was set at 4 Hz, i.e. in a cycle of 60 seconds, a total of $60\text{s} \times 4\frac{1}{\text{s}} = 240$ pairs of slopes and averages are recorded, which totals to 480 features per cycle. In other words, during one experiment – 4 cycles of 60 seconds – a total of $480 \times 4 = 1920$ features are extracted.

One way to visualize this process is shown in Figures 2.3. Note that the y-axis is in log-scale due to the different orders of magnitude of frequencies. Moreover, Figure 2.4 gives more insight into feature measurement and Table 2.1 summarizes the data acquisition details.

Table 2.1: Data acquisition details

Parameter	Value
Factors (gases)	3
Levels (concentrations)	5
Frequencies	16
Features per cycle	480
Number of cycles	4
Data points per mixture	1920
Number of mixtures	125
Features per experiment	240.000
Number of experiments	3
Total features	720.000

For specific timestamps and measurement durations, the reader is referred to Appendix A.

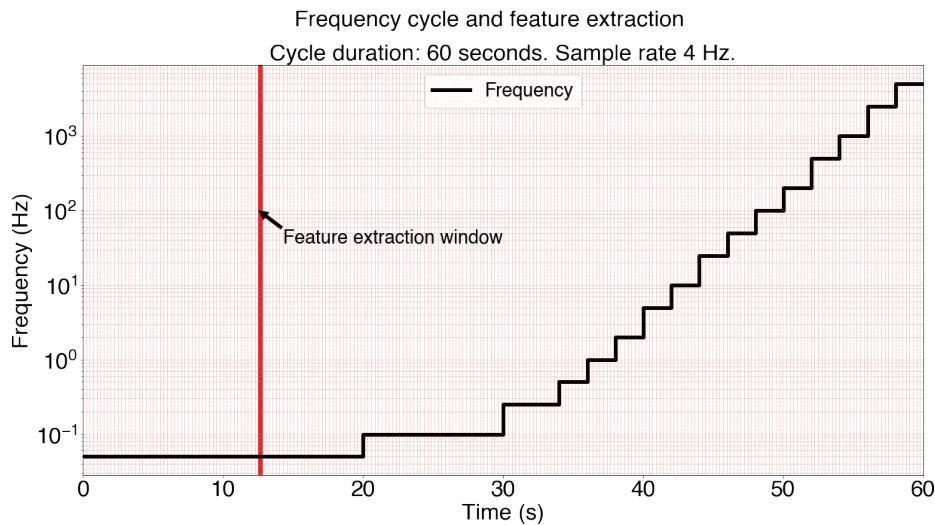


Figure 2.3: Feature measurements times per cycle. The width of the red line indicates the duration of one of the feature measurement windows as an example.

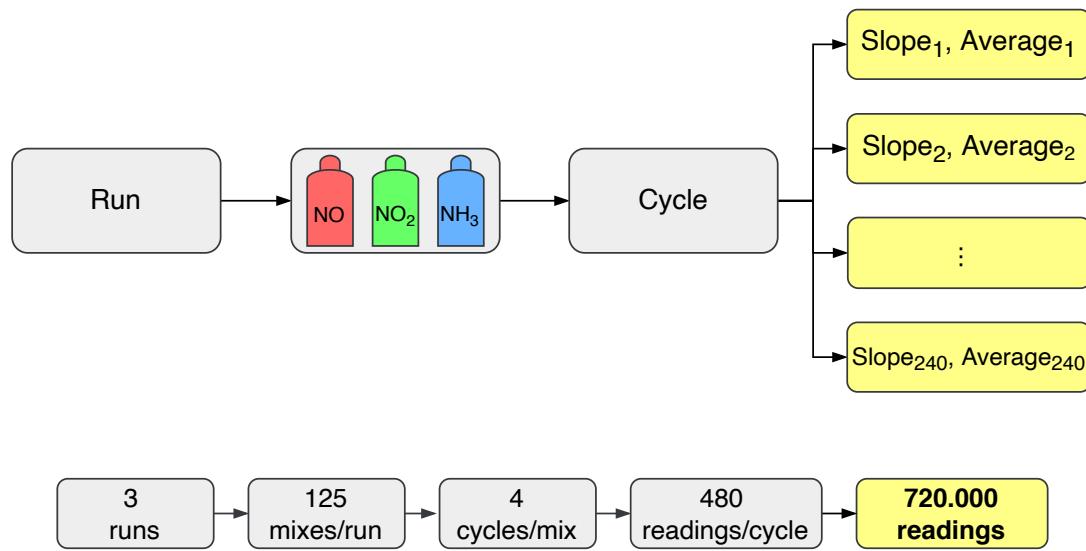


Figure 2.4: A visualization of the feature measurement process.

2.2 Raw data

The experiments were run between 26th and 29th March, 2021. The experiment data was exported as an excel file containing twelve columns, as specified in Table 2.2

Table 2.2: Raw data column details

Name	Description	Unit
Exposure nr	A particular mix of NO, NO ₂ and NH ₃ . Ranges from 1 to 375	-
Cycle nr	The cycle number. Ranges from 1 to 4.	-
Sample nr	Extracted feature index. Ranges from 1 to 240	-
NO	Nitric Oxide concentration	ppm
NO ₂	Nitrogen Dioxide concentration	ppm
NH ₃	Ammonia concentration	ppm
Freq	Frequency	Hz
Slope sensor 1	Slope	μA/s
Slope sensor 2	Slope	μA/s
Average sensor 1	Average	μA
Average sensor 2	Average	μA
Sensor temperature	Temperature	degrees Celsius (°C)

Table 2.3: Sample of raw data.

Index	Exposure nr	Cycle nr	Sample nr	NO [ppm]	NO ₂ [ppm]	NH ₃ [ppm]	Freq [Hz]	Slope sensor 1 [μA/s]	Slope sensor 2 [μA/s]	Average sensor 1 [μA/s]	Average sensor 2 [μA/s]	Sensor temperature [C]
0	1	1	1	10	5	20	0.05	-18.855169	-22.588416	32.926184	27.961554	274.994683
1	1	1	2	10	5	20	0.05	-28.289268	-28.185027	25.853867	20.915297	274.980487
2	1	1	3	10	5	20	0.05	-0.390916	-0.482129	25.756138	20.794765	274.985895
3	1	1	4	10	5	20	0.05	-0.234549	-0.156366	25.697501	20.755673	275.020372
4	1	1	5	10	5	20	0.05	-0.143336	-0.247580	25.661667	20.693778	275.014964
⋮												
100000	105	1	161	5	5	40	5.0	-38.366212	-48.495271	30.241896	24.821197	275.021724
100001	105	1	162	5	5	40	5.0	6.619507	8.521964	31.896773	26.951688	274.999415
100002	105	1	163	5	5	40	5.0	-1.941549	6.580416	31.411386	28.596792	275.011584
100003	105	1	164	5	5	40	5.0	27.401023	22.012900	38.261641	34.100017	275.009894
100004	105	1	165	5	5	40	5.0	-27.016623	-28.439121	31.507486	26.990236	275.014400
⋮												
359995	375	4	236	20	80	5	5000.0	-0.136821	-0.158538	34.129879	30.345597	275.002007
359996	375	4	237	20	80	5	5000.0	0.010859	0.010859	34.132593	30.348312	274.986797
359997	375	4	238	20	80	5	5000.0	-0.043435	0.030405	34.121734	30.355913	274.979811
359998	375	4	239	20	80	5	5000.0	-0.117275	-0.026061	34.092416	30.349398	274.984543
359999	375	4	240	20	80	5	5000.0	0.073840	0.039092	34.110876	30.359171	274.998063

2.3 Pre-processing

The features (slopes and averages) from the same target (a particular exposure) in the raw data file in Table 2.3 are spread across multiple rows, which is not suitable for analysis. Opposed to TCO, the experiments were conducted at constant temperature and therefore the temperature column is discarded. The data was, subsequently, modified to have the desired format: each row containing the predictors for one particular combination of gases. Additionally, the data from each sensor was split into two datasets.

The naming convention for the features is shown in Figure 2.5. First, the frequency in which the measurement was taken followed by the sensor number. After that, the feature name itself is followed by its index, i.e. where in the frequency cycle the measurement was made. This convention allows for easy identification of key information of the cycle and measurement.

The pre-processing results in the format shown in Figure 2.6. Recalling that there are 125 possible mixtures of gases, it is important to note that there are repeated exposures in the data set and those are treated as individual observations. Since each unique gas mixture was exposed 4 times during a cycle, and the experiment was repeated 3 times, this yields a total of $4 \times 3 \times 125 = 1500$ exposures.

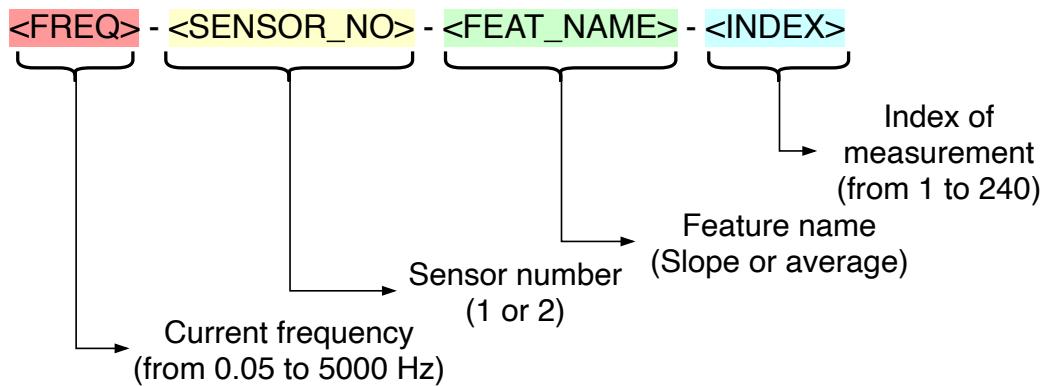


Figure 2.5: Feature naming convention.

Sensor 1			Sensor 2		
1500 exposures					
NO NO ₂ NH ₃			NO NO ₂ NH ₃		
Gases			Gases		
Slopes			Slopes		
Averages			Averages		
480 features					

Figure 2.6: Pre-processed data structure.

Finally, the majority of models proposed are scale variant, with the exception of OLS. The data, therefore, was centered and scaled to have unit variance and zero mean. A snippet of the final data set is shown in Table 2.4.

Table 2.4: Sample of pre-processed data.

Index	exposure	NO	NO2	NH3	0.05-1-slope-0	0.05-1-slope-1	...	5000.0-1-slope-239	0.05-1-avg-0	0.05-1-avg-1	0.05-1-avg-2	...	5000.0-1-avg-238	5000.0-1-avg-239
0	1.0	10.0	5.0	20.0	2.033884	-2.024181	...	1.256538	2.619349	1.757883	1.728610	...	2.276652	2.301890
1	1.0	10.0	5.0	20.0	-0.661288	0.869532	...	-0.846043	0.567699	2.477401	2.449619	...	2.522155	2.496103
2	1.0	10.0	5.0	20.0	0.232116	0.038766	...	1.546549	1.389922	2.716136	2.746701	...	2.637892	2.669321
3	1.0	10.0	5.0	20.0	-0.854716	1.071274	...	1.546549	0.592878	2.848766	2.866869	...	2.722064	2.753305
4	2.0	20.0	40.0	40.0	1.876790	-1.808299	...	-1.063551	2.963172	2.760346	2.746701	...	2.213523	2.182912
⋮														
700	176.0	40.0	20.0	40.0	1.735014	-1.850717	...	-2.151093	1.590775	0.071822	0.072959	...	0.098694	0.046562
701	176.0	40.0	20.0	40.0	-0.466435	0.439147	...	-1.208557	-0.393152	-0.046993	-0.033022	...	0.084665	0.055310
702	176.0	40.0	20.0	40.0	-0.509181	0.495187	...	-0.725205	-0.423540	-0.015217	-0.020505	...	0.078820	0.061143
703	176.0	40.0	20.0	40.0	-0.494932	0.467598	...	-2.296098	-0.409070	-0.031795	-0.043871	...	0.070052	0.014485
704	177.0	80.0	40.0	40.0	1.763512	-1.868305	...	2.005733	1.608140	0.076796	0.059607	...	0.074144	0.122381
⋮														
1495	374.0	80.0	80.0	40.0	-0.429032	0.485703	...	-0.725205	-1.122473	-1.365005	-1.369615	...	-1.476030	-1.490233
1496	375.0	20.0	80.0	5.0	1.083129	-1.205831	...	-1.426065	0.366557	-1.232375	-1.238876	...	-1.479537	-1.510646
1497	375.0	20.0	80.0	5.0	-0.635640	0.696067	...	0.507342	-1.303357	-1.373294	-1.352369	...	-1.461417	-1.445908
1498	375.0	20.0	80.0	5.0	-0.016883	0.086357	...	0.120661	-0.768521	-1.329084	-1.358488	...	-1.481875	-1.475653
1499	375.0	20.0	80.0	5.0	-0.588619	0.638821	...	0.676515	-1.247789	-1.358926	-1.369615	...	-1.486551	-1.466904



3 Theory

The quantification of gases based on the sensor response can be viewed as a multivariate multiple regression problem where the predictors, i.e. features derived from the sensor signal, are used to predict multiple responses, i.e. the concentrations of pertinent gases. This chapter exposes the theory behind some of these models.

The models here listed were chosen as a natural progression from a statisticians point of view: starting with simple models and progressively increasing complexity as insights from the data and the problem are gathered.

3.1 Notation

In favor of consistency and clarity, the notation used throughout this work is presented here. Bold capital letters, e.g. \mathbf{A} , are matrices while bold lower case letters are row vectors, e.g. \mathbf{a} . Scalars, on the other hand, are denoted as standard lower case letters, e.g. a_{i1} . Transposes and inverses are denoted, respectively with \cdot^\top and \cdot^{-1} . An example is shown below.

$$\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]^\top$$

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_n \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ a_{21} & a_{22} & \dots & a_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{np} \end{bmatrix}$$

3.2 Ordinary Least Squares Regression

A simple, first approach would be to tackle the problem with a Ordinary Least Squares (OLS) regression model. As Friedman, Hastie, Tibshirani, et al. 2001 explains, each output $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K]^T$ has its own linear model. Now, given a set of n observations $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$ and each observation having p features, e.g. $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{ip}]$, $i = 1, 2, \dots, n$, the concatenation of all linear models can be written in matrix form as in Equation 3.1.

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E} \quad (3.1)$$

Where:

- \mathbf{B} : $[p+1 \times K]$ matrix of regression coefficients (with the $+1$ referring to the intercept term);
- \mathbf{E} : $[N \times K]$ matrix of residuals.

The objective is then to find the coefficients $\hat{\mathbf{B}}$ which minimizes the Residual Sum of Squares (RSS), which is summarized by Equation 3.3 (Friedman, Hastie, Tibshirani, et al. 2001):

$$\hat{\mathbf{B}}^{\text{OLS}} = \arg \min_{\mathbf{B}} \text{RSS}(\mathbf{B}) \quad (3.2)$$

In turn, the RSS, as the name suggests, is defined as the difference between real and predicted values, squared, which in matrix form is written as (Friedman, Hastie, Tibshirani, et al. 2001):

$$\text{RSS}(\mathbf{B}) = \text{Tr}[(\mathbf{Y} - \mathbf{XB})^T(\mathbf{Y} - \mathbf{XB})] \quad (3.3)$$

Finally, solving for $\hat{\mathbf{B}}$ yields (Friedman, Hastie, Tibshirani, et al. 2001):

$$\hat{\mathbf{B}}^{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (3.4)$$

For the problem at hand, in addition to the high number of features, it is often the case that sensor data points are acquired in quick succession, which in turn leads to highly correlated features (Bastuck 2019), which can result in high variance in a least squares model (Friedman, Hastie, Tibshirani, et al. 2001). It is natural, therefore, to progress towards methods that incorporate dimensionality reduction such as Principal Components Regression (PCR) and Partial Least Squares Regression (PLSR) or shrinkage such as Ridge Regression.

3.3 Principal Component Analysis

One way to define Principal Components Analysis (PCA) is to view it as a orthogonal projection of the data into a principal space of lower dimension such that the variance of this projection is maximized (Bishop 2006).

Just as before, consider the collection of n observations \mathbf{X} with covariance matrix Σ . Additionally, consider a matrix $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n]^\top$ where \mathbf{p}_i is a row vector of coefficients referring to the i -th linear combination (Johnson and Wichern 2013):

$$t_i = \mathbf{x}_i \mathbf{p}_i^\top \quad i = 1, 2, \dots, n \quad (3.5)$$

The variance and covariance of these new variables t_i can be written as follows:

$$\text{Var}(t_i) = \mathbf{p}_i^\top \Sigma \mathbf{p}_i \quad i = 1, 2, \dots, n \quad (3.6)$$

$$\text{Cov}(t_i, t_k) = \mathbf{p}_i^\top \Sigma \mathbf{p}_k \quad i, k = 1, 2, \dots, n \quad (3.7)$$

The first Principal Component (PC) is then the linear combination with maximum variance, i.e. the linear combinations that maximizes $\text{Var}(t_1)$, with the constraint that the coefficient vector \mathbf{p}_1 has unit length. In summary, the first PC is computed as (Johnson and Wichern 2013):

$$\begin{aligned} t_1 &= \mathbf{x}_1 \mathbf{p}_1^\top \\ &\text{that maximizes } \text{Var}(\mathbf{x}_1 \mathbf{p}_1^\top) \\ &\text{subject to } \mathbf{p}_1^\top \mathbf{p}_1 = 1 \end{aligned} \quad (3.8)$$

The second PC, similarly to the first, is the linear combination with maximum variance, but with an added extra constraint: this new linear combination must be orthogonal to the previous one, i.e. they must be linearly independent:

$$\begin{aligned} t_2 &= \mathbf{x}_2 \mathbf{p}_2^\top \\ &\text{that maximizes } \text{Var}(\mathbf{x}_2 \mathbf{p}_2^\top) \\ &\text{subject to } \mathbf{p}_2^\top \mathbf{p}_2 = 1 \\ &\text{and } \text{Cov}(t_1, t_2) = 0 \end{aligned} \quad (3.9)$$

The k -th PC is then:

$$\begin{aligned} t_k &= \mathbf{x}_k \mathbf{p}_k^\top \\ &\text{that maximizes } \text{Var}(\mathbf{x}_k \mathbf{p}_k^\top) \\ &\text{subject to } \mathbf{p}_k^\top \mathbf{p}_k = 1 \\ &\text{and } \text{Cov}(t_j, t_k) = 0 \text{ for } k > j \end{aligned} \quad (3.10)$$

In summary, the objective of PCA is find a matrix \mathbf{P} such that the linear transformation

$$\mathbf{T} = \mathbf{X} \mathbf{P}^\top \quad (3.11)$$

yields new variables that are uncorrelated and arranged in decreasing order of variance.

It can be shown that these desired linear combinations can be written in terms of the eigenvalues (λ) and eigenvectors (e) of Σ , the covariance matrix of X (Johnson and Wichern 2013). The elements of eigenvectors are called loadings, while the new features T are . In short, for the k -th PC:

$$\begin{aligned} t_k &= X_k e_k^\top \\ \text{Var}(t_k) &= e_k^\top \Sigma e_k = \lambda_k \\ \text{Cov}(t_j, t_k) &= e_k^\top \Sigma e_j = 0 \quad \text{for } k \neq j \end{aligned} \tag{3.12}$$

There are several ways of computing PCs. Many of which involving finding aforementioned eigenvalues and eigenvectors. These calculations can be computationally expensive, depending on the desired number of extracted PCs (Bishop 2006). One option is the Nonlinear Iterative Partial Least Squares (NIPALS) algorithm, also called Power Method. It has two clear advantages: "it can handle missing data and computes the components sequentially" (Dunn 2021).

The NIPALS algorithm to compute the first k -th PCs, PC_i , $i = 1, 2, \dots, k$ us displayed below as Algorith 1 (Dunn 2021) (Ng 2013) (Wright 2017). Since it computes the loadings and scores sequentially, it is possible to stop it as early as desired. The "truncated" loadings and scores that project X into the principal subspace of k PCs is defined in Equation 3.13 :

$$T_{|k} = X P_{|k}^\top \tag{3.13}$$

Algorithm 1: Nonlinear Iterative Partial Least Squares (NIPALS) for PCA

Result: Matrices of loadings $P_{|k}$ and scores $T_{|k}$ of the k -th first Principal Components

```

1 Initialize T_{|k} and P_{|k}
2 i = 1
3 X_1 := X
4 while i < k do
5   repeat
6     Choose t_i as any column of X_i
7     Compute loadings p_i = (t_i^\top t_i)^{-1} t_i^\top X_i
8     Scale p_i = p_i / sqrt(p_i^\top p_i)
9     Compute scores t_i = (p_i^\top p_i)^{-1} p_i^\top X_i
10  until t_i converges
11  Append t_i to T_{|k}
12  Append p_i to P_{|k}
13  Deflate: X_{i+1} = X_i - t_i p_i^\top
14  i += 1
15 end
16 return T_{|k}, P_{|k}

```

In words, the main idea of the algorithm can be summarized as choosing an arbitrary column of \mathbf{X} as the scores vector t_i , shown in line 6. Next, the computation of the i -th loadings vector p_i by regressing every column of \mathbf{X} via OLS onto the scores t_i . p_i is then scaled to have unit length in Line 8, which in turn is used to compute the i -th scores vector t_i by regressing every column of \mathbf{X} via OLS onto the loadings p_i , shown in Line 9. This procedure is repeated until change in t_i between iterations is small enough. Once convergence is achieved, scores t_i and loadings p_i are stored as the i -th column of matrices \mathbf{T} and \mathbf{P} of Equation 3.11, respectively. Finally, the variability explained by t_i and p_i from \mathbf{X} is subtracted in a procedure called deflation.

3.4 Principal Component Regression

With the inner workings of PCA explained in the previous section, PCR can be simply reduced to a Least Squares regression on the first k -th PCs, i.e. performing linear regression on $\mathbf{T}_{|k}$ instead of \mathbf{X} :

$$\mathbf{Y} = \mathbf{T}_{|k}\mathbf{B} + \mathbf{E} \quad (3.14)$$

And the regression coefficients are found analogously to Equation 3.4:

$$\hat{\mathbf{B}}^{\text{PCR}} = (\mathbf{T}_{|k}^\top \mathbf{T}_{|k})^{-1} \mathbf{T}_{|k}^\top \mathbf{Y} \quad (3.15)$$

Although useful, PCR has a potential flaw: while the new found projection of \mathbf{X} is guaranteed to best explain the variance of predictors, this cannot be said about the responses \mathbf{Y} (James, Witten, Hastie, and Tibshirani 2013). PLSR, on the other hand, solves this issue by supervising the identification of PCs (James, Witten, Hastie, and Tibshirani 2013).

3.5 Partial Least Squares Regression

PLSR, much like PCR, also aims to reduce dimensionality via linear combinations of the inputs. This technique, however, also takes into account the response variables \mathbf{Y} . One key advantage of PLSR is that it seeks axes with most variance (like PCR) and high correlation with response variables (Friedman, Hastie, Tibshirani, et al. 2001).

The main idea can be described as finding linear combinations for the data matrix \mathbf{X} and response matrix \mathbf{Y} as follows (Ng 2013), similarly to what was done in Section 3.3. Here, the matrices \mathbf{W} and \mathbf{U} are score matrices, i.e. the transformed PLS variables, and \mathbf{L} and \mathbf{Q} are loading matrices, i.e. the weights of this transformation (projection).

$$\mathbf{W} = \mathbf{XL}^\top \quad (3.16)$$

$$\mathbf{U} = \mathbf{YQ}^\top \quad (3.17)$$

Instead of simply running NIPALS on \mathbf{X} and \mathbf{Y} separately, PLSR uses information from \mathbf{Y} to decompose \mathbf{X} and vice-versa (Ng 2013). Algorithm 2 is an adaptation of Algorithm 1 to incorporate this intended behavior.

Algorithm 2: NIPALS for Partial Least Squares Regression (PLSR)

Result: Matrices of loadings $\mathbf{L}_{|k}$, $\mathbf{Q}_{|k}$ and scores $\mathbf{W}_{|k}$, $\mathbf{U}_{|k}$ of the k -th first Partial Least Squares directions

```

1 Initialize  $\mathbf{L}_{|k}$ ,  $\mathbf{Q}_{|k}$  and  $\mathbf{W}_{|k}$ ,  $\mathbf{U}_{|k}$ 
2  $i = 1$ 
3  $\mathbf{X}_1 := \mathbf{X}$ 
4  $\mathbf{Y}_1 := \mathbf{Y}$ 
5 while  $i < k$  do
6   repeat
7     Choose  $\mathbf{u}_i$  as any column of  $\mathbf{Y}_i$ 
8     Compute loadings of  $\mathbf{X}_i$  based on score of  $\mathbf{Y}_i$ :  $\ell_i = (\mathbf{u}_i^\top \mathbf{u}_i)^{-1} \mathbf{u}_i^\top \mathbf{X}_i$ 
9     Scale  $\ell_i = \frac{\ell_i}{\sqrt{\ell_i^\top \ell_i}}$ 
10    Compute score of  $\mathbf{X}_i$ :  $\mathbf{w}_i = (\ell_i^\top \ell_i)^{-1} \ell_i^\top \mathbf{X}_i$ 
11    Compute loadings of  $\mathbf{Y}_i$  based on score of  $\mathbf{X}_i$ :  $\mathbf{q}_i = (\mathbf{w}_i^\top \mathbf{w}_i)^{-1} \mathbf{w}_i^\top \mathbf{Y}_i$ 
12    Scale  $\mathbf{q}_i = \frac{\mathbf{q}_i}{\sqrt{\mathbf{q}_i^\top \mathbf{q}_i}}$ 
13    Compute score of  $\mathbf{Y}_i$ :  $\mathbf{u}_i = (\mathbf{q}_i^\top \mathbf{q}_i)^{-1} \mathbf{q}_i^\top \mathbf{Y}_i$ 
14  until  $\mathbf{u}_i$  converges
15  Append  $\mathbf{w}_i$  to  $\mathbf{W}_{|k}$ 
16  Append  $\ell_i$  to  $\mathbf{L}_{|k}$ 
17  Append  $\mathbf{u}_i$  to  $\mathbf{U}_{|k}$ 
18  Append  $\mathbf{q}_i$  to  $\mathbf{Q}_{|k}$ 
19  Deflate  $\mathbf{X}_i$ :  $\mathbf{X}_{i+1} = \mathbf{X}_i - \mathbf{w}_i \ell_i^\top$ 
20  Deflate  $\mathbf{Y}_i$ :  $\mathbf{Y}_{i+1} = \mathbf{Y}_i - \mathbf{u}_i \mathbf{q}_i^\top$ 
21   $i += 1$ 
22 end
23 return  $\mathbf{W}_{|k}$ ,  $\mathbf{L}_{|k}$ ,  $\mathbf{U}_{|k}$ ,  $\mathbf{Q}_{|k}$ 

```

As with Algorithm 1, Algorithm 2 can be summarized as choosing a column of \mathbf{Y}_i as the initial response score vector \mathbf{u}_i . After that, the i -th loadings vector \mathbf{w}_i of \mathbf{X} is computed in Line 8 by regressing every column of \mathbf{X} via OLS onto scores vector of \mathbf{Y} , \mathbf{u}_i . Similarly to before, the data loadings vector \mathbf{w}_i is scaled to have unit length, which in turn is used to compute the i -th data scores vector \mathbf{w}_i by regressing every column of \mathbf{X}_i via OLS onto the column ℓ_i in Line 10. Now, the i -th response loadings vector \mathbf{q}_i of \mathbf{Y}_i by regressing every column of \mathbf{Y} via OLS onto scores vector of \mathbf{X} , \mathbf{w}_i , shown in Line 11; This loadings vector is also scaled to have unit length.

Following in Line 13, the i -th response scores vector \mathbf{u}_i is computed by regressing every column of \mathbf{Y}_i via OLS onto the column \mathbf{q}_i . This procedure is repeated until change in \mathbf{u}_i between iterations is small enough. In that case the results \mathbf{w}_i and ℓ_i are stored as the i -th

column of matrices \mathbf{W} and \mathbf{L} of Equation 3.16 and \mathbf{u}_i and \mathbf{q}_i are stored as the i -th column of matrices \mathbf{U} and \mathbf{Q} of Equation 3.17. Finally, the variability explained by \mathbf{w}_i, ℓ_i and $\mathbf{u}_i, \mathbf{q}_i$ from \mathbf{X}_i and \mathbf{Y}_i , respectively, are removed.

After finding the k partial least squares directions from Algorithm 2 above, the problem, as in Section 3.4, reduces to performing Least Squares Regression using the newfound transformations.

$$\mathbf{Y} = \mathbf{W}_{|k}\mathbf{B} + \mathbf{E} \quad (3.18)$$

Which in turn, analogously to Equations 3.4 and 3.15, yields the coefficients:

$$\hat{\mathbf{B}}^{\text{PLSR}} = (\mathbf{W}_{|k}^\top \mathbf{W}_{|k})^{-1} \mathbf{W}_{|k}^\top \mathbf{Y} \quad (3.19)$$

3.6 Ridge Regression

Ridge regression is also a viable alternative to reduce the problem of highly correlated features (Friedman, Hastie, Tibshirani, et al. 2001). Instead of fitting a least squares model on a subset of predictors or a transformation of them, Ridge allows the use of all features with a continuous shrinkage of its coefficients, which results in less variance (Friedman, Hastie, Tibshirani, et al. 2001).

For the multi-output case, there are two options: use the same penalization parameter λ for all variables $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K]^\top$ or apply different parameters $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \dots, \lambda_K]^\top$. In this work, the latter is preferred over the former, as it allows a more fine tuned control of the regression models for each studied gas.

Analogous to Section 3.2, the goal is to minimize the RSS, but now with the penalization term taken into account. Equation 3.20 below shows this objective function in matrix form.

$$\text{RSS}^{\text{Ridge}}(\mathbf{B}, \boldsymbol{\lambda}) = \text{Tr}[(\mathbf{Y} - \mathbf{XB})^\top (\mathbf{Y} - \mathbf{XB})] + \text{Tr}[\mathbf{B}^\top \mathbf{B} + \boldsymbol{\lambda} \mathbf{I}] \quad (3.20)$$

$$\hat{\mathbf{B}}^{\text{Ridge}} = \arg \min_{\mathbf{B}} \text{RSS}^{\text{Ridge}}(\mathbf{B}) \quad (3.21)$$

The coefficients that minimize the RSS is shown in Equation ?? below.

$$\hat{\mathbf{B}}^{\text{Ridge}} = (\mathbf{X}^\top \mathbf{X} + \boldsymbol{\lambda} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y} \quad (3.22)$$

It is important to note that the parameters β_0 are calculated separately as the response mean \bar{Y} (Friedman, Hastie, Tibshirani, et al. 2001), i.e.:

$$\hat{\mathbf{B}}_0^{\text{Ridge}} \begin{bmatrix} \beta_{0,1} \\ \beta_{0,2} \\ \vdots \\ \beta_{0,K} \end{bmatrix} = \begin{bmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \vdots \\ \bar{y}_K \end{bmatrix} = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n y_{i,1} \\ \frac{1}{n} \sum_{i=1}^n y_{i,2} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n y_{i,K} \end{bmatrix} \quad (3.23)$$

The choice of hyper-parameters $\lambda \geq 0$ controls how much shrinkage is applied to the coefficients: larger λ implies more penalization to complex models. Although the coefficients are shrunk towards zero, they never reach zero, which makes Ridge regularization unsuitable for feature selection (Friedman, Hastie, Tibshirani, et al. 2001).

Finally, predictions in a ridge regression setting are computed as:

$$\hat{\mathbf{Y}}^{\text{Ridge}} = \hat{\mathbf{B}}_0^{\text{Ridge}} + \hat{\mathbf{B}}^{\text{Ridge}} \mathbf{X} \quad (3.24)$$

3.7 Cross Validation

There are several choices to make for the aforementioned models: How many PCs/PLS components to use? How much penalization to impose in Ridge regression?

A first answer to this would be to split the data into training and validation sets. After fitting models to the former, the latter is used to measure the prediction error via some scoring function. In that sense, it is important to distinguish test error rate from training error rate. The first, also called generalization error, is the score of the fit on an independent, previously unseen test sample. The second, on the other hand, is the average score over the training sample (Friedman, Hastie, Tibshirani, et al. 2001).

Scoring functions measure how much the data deviates from the fit and can be used as a qualitative tool for model selection and comparison. Once this is done, the choice of the model that yields minimum error is trivial. Two examples of wildly used score functions are Mean Squared Error (MSE) and Root Mean Squared Error (RMSE). For the multi-output case of m responses and n observations, they are defined respectively as :

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^m (y_{ij} - \hat{y}_{ij}) \right)^2 \quad (3.25)$$

$$\text{RMSE} = \sqrt{\text{MSE}} \quad (3.26)$$

This approach, however, is sensitive to the choice of these sets. Additionally, reserving part of the data just for validation might be detrimental for the model fitting process, specially if the number of observations is low (James, Witten, Hastie, and Tibshirani 2013).

One tool that can help alleviating these problems is Cross Validation (CV). More specifically, K-fold CV: it involves equally dividing the training data into K sets. For each subset, the

desired model is trained using $k - 1$ folds and the prediction error is computed on the remaining fold (Friedman, Hastie, Tibshirani, et al. 2001). As for the final evaluation, it is performed in the held-out test set.



4 Methods

This chapter describes the application of theory shown in Chapter 3 applied to the pre-processed data described in Chapter 2. For further details regarding coding and implementation, the reader is referred to this work's [repository](#).

In general, the use of Scikit-Learn's pipeline class alongside linear models made analysis straightforward. Additionally, Scikit-Learn's GridSearchCV allowed for a faster evaluation of different hyperparameters such as number of components and shrinkage factor.

Throughout all methods, training, test and validation sets were split using `train_test_split`. Out of all 1500 exposures, 1200 (80%) were assigned to the training set and 300 (20%) to the validation set. The test set is embedded in the training set, as k-fold cross-validation is used. A fixed random seed of 42 was set to ensure reproducibility of results.

For reasons that will become clear in Chapter 6 - Discussions, regression analysis was done with two sets of predictor variables: first with all features, i.e. slopes and averages for each exposure, and second with only averages for unique mixtures (averaged features).

4.1 Ordinary Least Squares

Here treated as a baseline, OLS is fit using `LinearRegression` then evaluated using the validation set.

4.2 Principal Components Regression

First, a PCA was conducted via Scikit-Learn's `PCA`. With that, cumulative variance and score plots were made in an attempt of better visualizing and understanding the data.

Following that, a linear regression on the PCs was made. The number of components to use in it was found via CV with RMSE as the scoring function. After training, the model was evaluated in a held-out validation set and a actual vs. predicted plot was constructed.

4.3 Partial Least Squares Regression

Here, a similar procedure to Section 4.2 was conducted. Initially, PLS components were extracted via Scikit-learn's `PLSRegression` and some informative plots were made: cumulative explained variance (for both predictors \mathbf{X} and response \mathbf{Y}), as well as score plots for the first two PLS components in an attempt to better visualize the data. The regression model was trained with the ideal number of components given by CV and later evaluated in the held-out validation set.

4.4 Ridge Regression

Once again, cross-validation was used to select the amount of shrinkage, i.e. λ . After setting λ , analysis proceeded as usual: fitting the regression line via `Ridge` to training data and evaluating it using the validation set.

4.5 Averaging features

In efforts to further analyze the data, the previous 1500 observations are averaged by unique mixtures, i.e. for each mixture, the features are averaged from its twelve exposures, yielding 125 observations. Figure 4.1 clarifies this further.

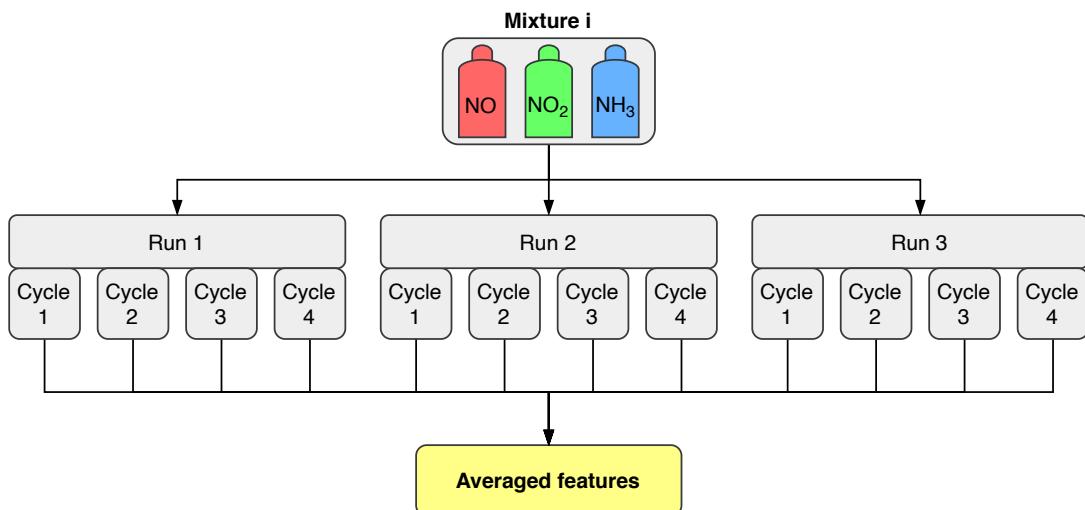


Figure 4.1: A visualization of the feature averaging process.



5 Results

This chapter is dedicated to show the analysis' results. In favor of clarity and organization, this chapter will be divided into two main sections: the first treating each gas exposure individually and the other averaging exposures of the same gas mixture. Subsections on the models used also are present. The plots presented in this section were made using *ad-hoc* plotting functions.

Initially the regression analysis was done in the pre-processed data presented in Chapter 2, i.e. each observation corresponds to a gas exposure. It is important to remind the reader that in this data, each unique gas mixture was exposed (i.e. an exposure) twelve times: four frequency cycles through three experiment repetitions, yielding 1500 observations.

Before beginning analysis, an assessment of correlation between features is first conducted and shown in Figure 5.1.

5.1 Individual exposures

5.1.1 OLS

As explained in Chapter 4, OLS is treated here as a baseline. The actual vs. predicted plot in Figure 5.2 shows the predictions of this model.

5.1.2 PCR

Following the methodology of Chapter 4, a PCA is conducted with two components in an attempt to visualize the data in a lower dimensional space in Figure 5.3.

Furthermore, an explained variance plot is shown in Figure 5.4. The first two PCs explain approximately 40% of the total variance, reaching 80% around 100 components.

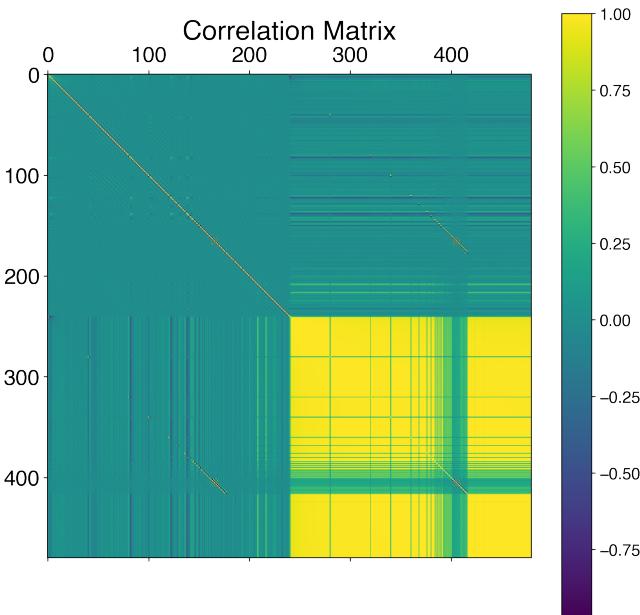


Figure 5.1: Correlation matrix of features.

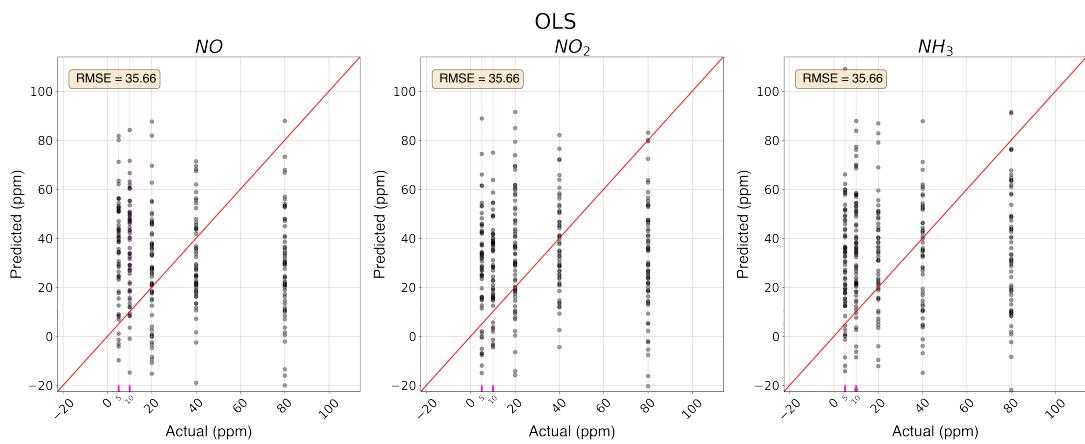


Figure 5.2: OLS prediction for individual exposures

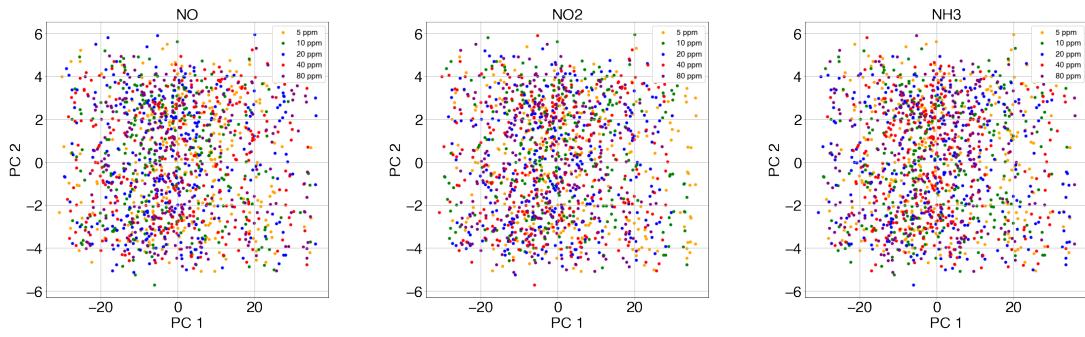


Figure 5.3: PCA with two components for the three gases.

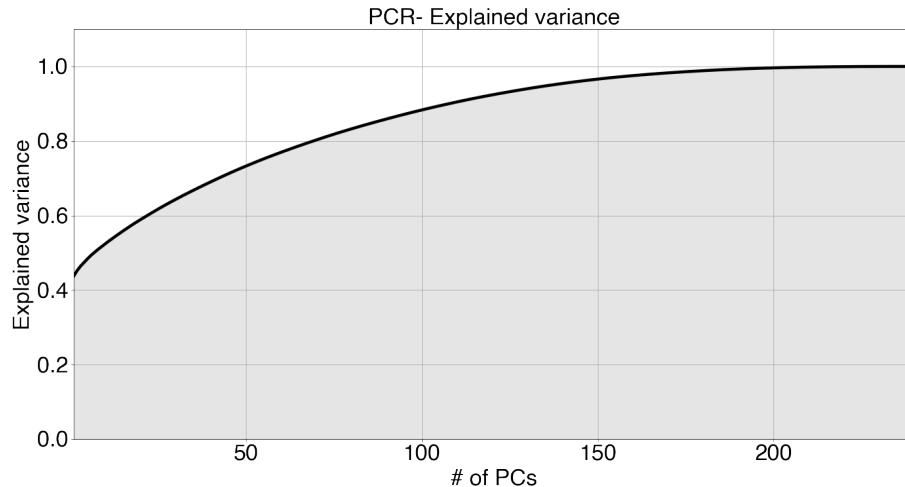


Figure 5.4: Explained variance

After this exploration of PCA, the analysis proceeds to fit a PCR model to the data. The choice of number of PCs was made via cross-validation using RMSE as the loss function, as can be seen in Figure 5.5. Choosing only one components yields the minimum loss, around 27.

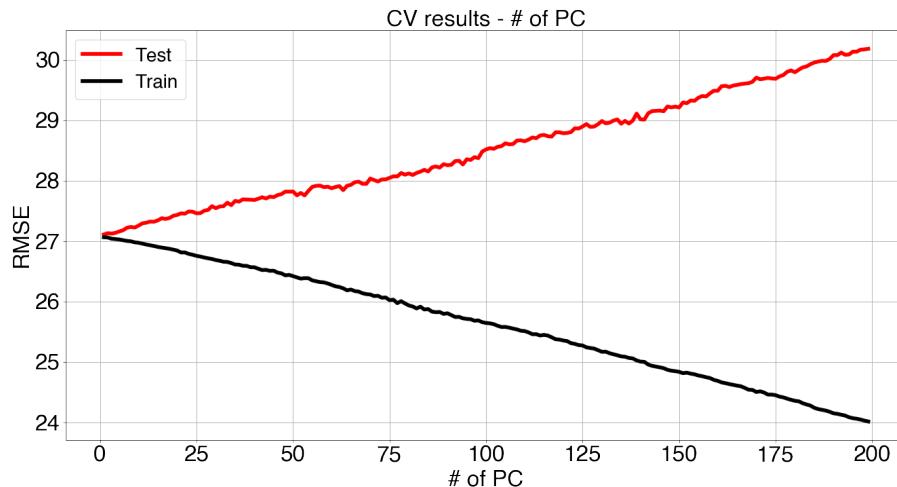


Figure 5.5: Number of PCs selection via CV.

After the choosing the number of components, the regression line was fit to the training data and used to predict unseen validation data. A quick way to visualize and assess the quality of the fit is an actual vs. predicted plot, shown in Figure 5.6.

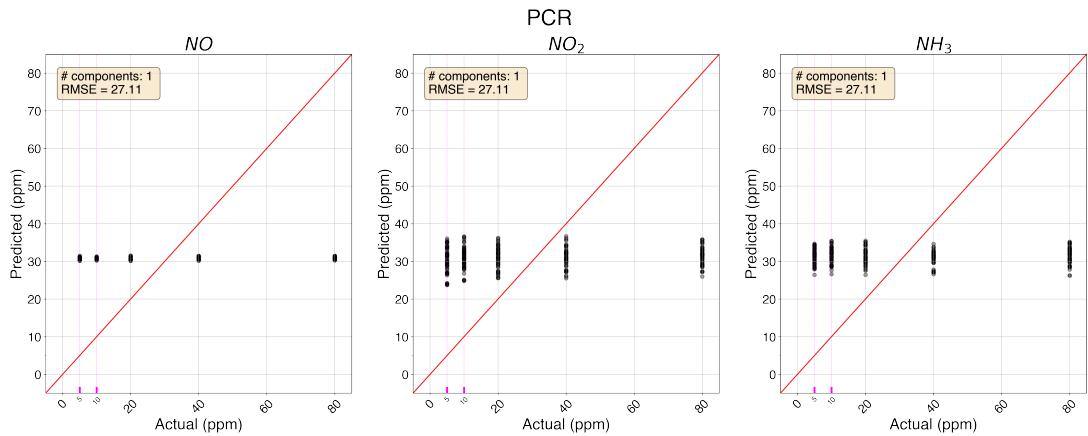


Figure 5.6: PCR prediction for individual exposures

5.1.3 PLSR

Following the proposed model progression, the analysis proceeds to fit the PLSR model. A similar pipeline to Section 5.1.2 was used. First, in Figure 5.7, the choice of only two PLS components allowed visualization of data in a two dimensional plot. Moreover, total explained variance is shown in Figure 5.8. Once again, cross-validation using RMSE yields a single component as the best choice, which is then used to fit and predict gas concentrations in Figure 5.10.

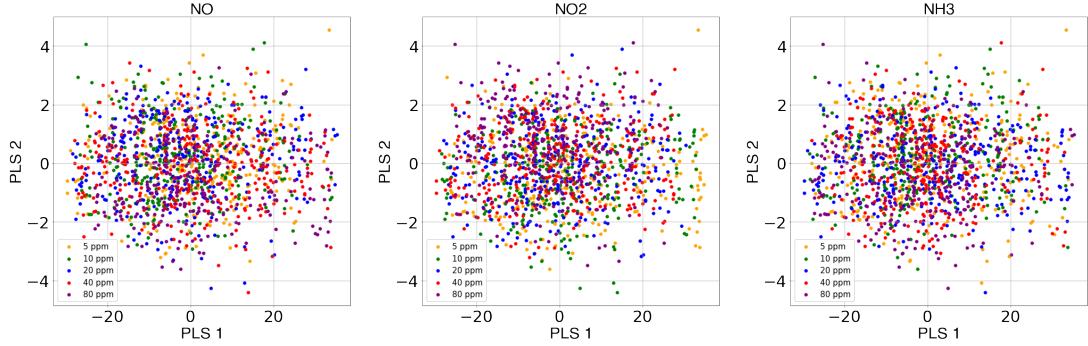


Figure 5.7: PLS with two components for the three gases.

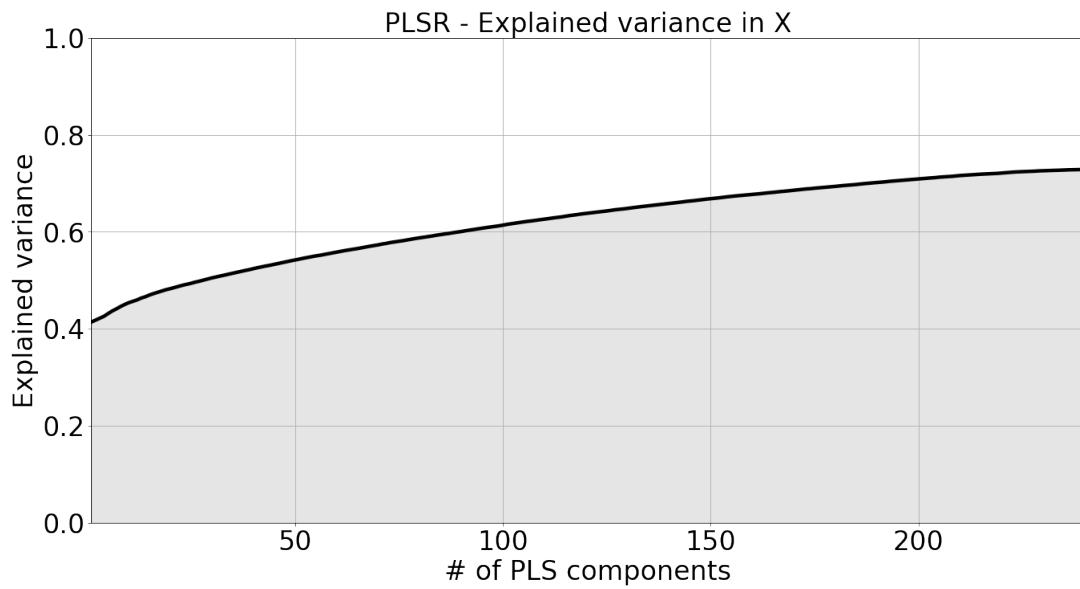


Figure 5.8: PLS - Explained variance of X

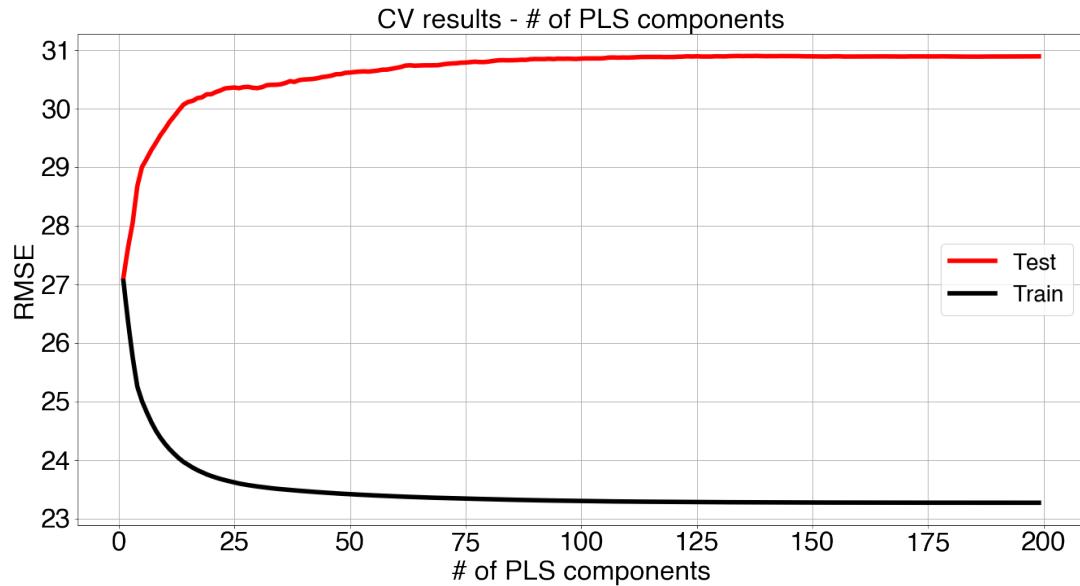


Figure 5.9: Number of PLS components selection via CV.

5.1.4 Ridge Regression

For Ridge regression, the regularization term λ was chosen via cross-validation, as shown in Figure 5.11. Additionally, the shrinkage of coefficients can be seen in Figure 5.12. As expected, the coefficients shrink asymptotically towards zero.

Finally, after the choice of λ , the actual vs. predicted plot is presented in Figure 5.13.

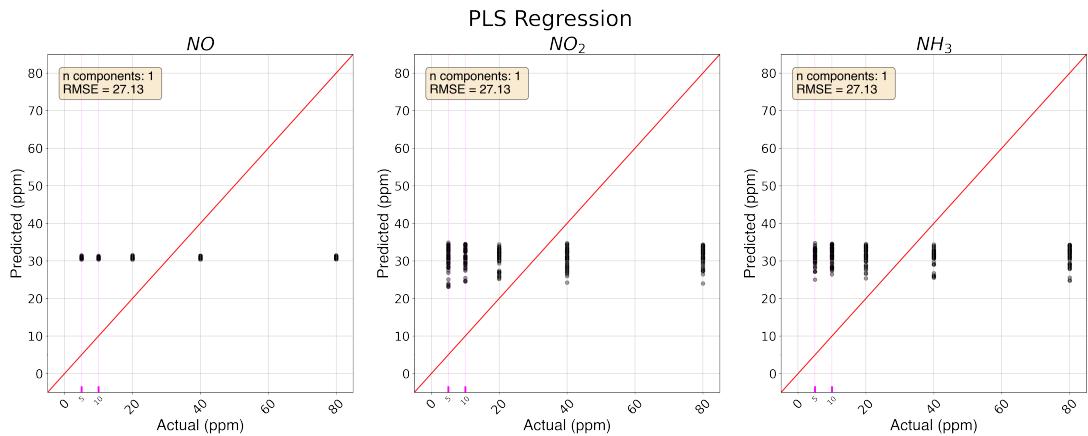
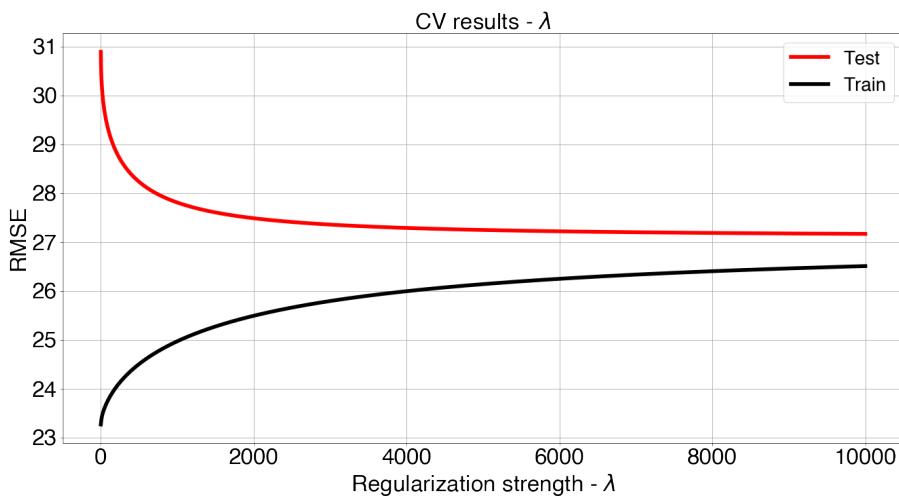


Figure 5.10: PLSR prediction for individual exposures


 Figure 5.11: CV of shrinkage factor λ

5.2 Averaging features by gas mixture

Following the process depicted in Figure 4.1, the averaged features data set is shown below on Table 5.1. Now it is now possible to plot slope and average features through time, as shown in Figures 5.14 and 5.15. Furthermore, for the following analysis, only average features were selected, i.e. slopes were discarded, which is shown in Table 5.2

Following that, all models were fit again to this new data set. The results are shown in the following subsections.

Table 5.1: Sample of averaged mixture data.

Index	NO	NO2	NH3	0.05-1-slope-0	0.05-1-slope-1	0.05-1-slope-2	...	5000.0-1-slope-239	0.05-1-avg-0	0.05-1-avg-1	...	5000.0-1-avg-238	5000.0-1-avg-239
0	5.0	5.0	5.0	-26.386178	-14.165356	-0.234640	...	-0.041897	28.983749	25.442410	...	35.162932	35.152458
1	5.0	5.0	10.0	-26.294332	-14.421532	-0.216090	...	0.004524	28.538652	24.933269	...	34.622460	34.623591
2	5.0	5.0	20.0	-25.514491	-15.174588	-0.253371	...	-0.010135	29.038925	25.245278	...	35.025637	35.023103
3	5.0	5.0	40.0	-25.906221	-14.565139	-0.306851	...	0.003801	28.698684	25.057399	...	34.575699	34.576649
4	5.0	5.0	80.0	-26.731849	-13.795072	-0.241065	...	-0.022803	28.738748	25.289980	...	34.860040	34.854340
:													
50	20.05.0	5.0	-25.541819	-15.120475	-0.228034	...	-0.004344	28.490467	24.710348	...	34.279956	34.278870	
51	20.05.0	10.0	-26.036798	-14.578261	-0.202154	...	-0.067958	28.623668	24.979102	...	34.562216	34.545227	
52	20.05.0	20.0	-25.982142	-14.508765	-0.221338	...	-0.026222	28.764243	25.137052	...	34.752584	34.746928	
53	20.05.0	40.0	-25.964768	-14.674361	-0.179532	...	0.007601	28.229947	24.561356	...	34.258736	34.260636	
54	20.05.0	80.0	-25.317133	-15.287429	-0.308570	...	-0.052484	29.162557	25.340699	...	35.015909	35.002788	
:													
120	80.080.0	5.0	-27.073901	-13.562242	-0.218533	...	0.028142	28.548244	25.157684	...	34.742313	34.749349	
121	80.080.0	10.0	-26.329623	-14.337196	-0.120442	...	-0.012669	28.630183	25.045884	...	34.678857	34.675690	
122	80.080.0	20.0	-25.935631	-14.734990	-0.196363	...	-0.027890	28.420835	24.737087	...	34.354338	34.347416	
123	80.080.0	40.0	-25.821070	-14.855702	-0.241336	...	-0.032395	28.457189	24.743263	...	34.327938	34.319839	
124	80.080.0	80.0	-26.503363	-14.087625	-0.186228	...	-0.038820	28.615161	25.093255	...	34.743829	34.734124	

Table 5.2: Sample of only average features.

Index	NO	NO2	NH3	0.05-1-avg-0	0.05-1-avg-1	...	5000.0-1-avg-238	5000.0-1-avg-239
0	5.0	5.0	5.0	28.983749	25.442410	...	35.162932	35.152458
1	5.0	5.0	10.0	28.538652	24.933269	...	34.622460	34.623591
2	5.0	5.0	20.0	29.038925	25.245278	...	35.025637	35.023103
3	5.0	5.0	40.0	28.698684	25.057399	...	34.575699	34.576649
4	5.0	5.0	80.0	28.738748	25.289980	...	34.860040	34.854340
:								
50	20.0	5.0	5.0	28.490467	24.710348	...	34.279956	34.278870
51	20.0	5.0	10.0	28.623668	24.979102	...	34.562216	34.545227
52	20.0	5.0	20.0	28.764243	25.137052	...	34.752584	34.746928
53	20.0	5.0	40.0	28.229947	24.561356	...	34.258736	34.260636
54	20.0	5.0	80.0	29.162557	25.340699	...	35.015909	35.002788
:								
120	80.0	80.0	5.0	28.548244	25.157684	...	34.742313	34.749349
121	80.0	80.0	10.0	28.630183	25.045884	...	34.678857	34.675690
122	80.0	80.0	20.0	28.420835	24.737087	...	34.354338	34.347416
123	80.0	80.0	40.0	28.457189	24.743263	...	34.327938	34.319839
124	80.0	80.0	80.0	28.615161	25.093255	...	34.743829	34.734124

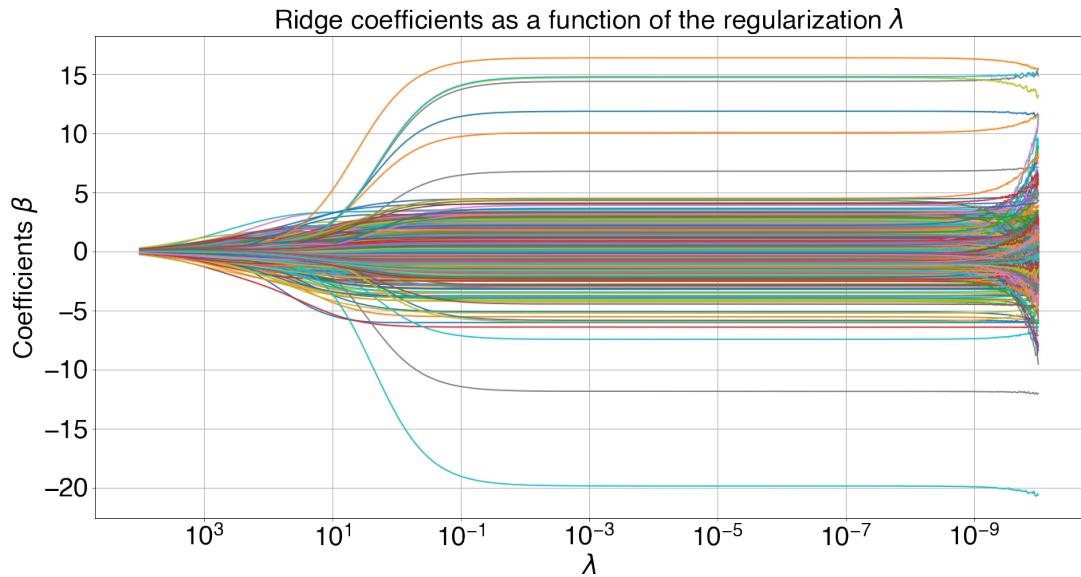


Figure 5.12: Shrinkage of ridge coefficients.

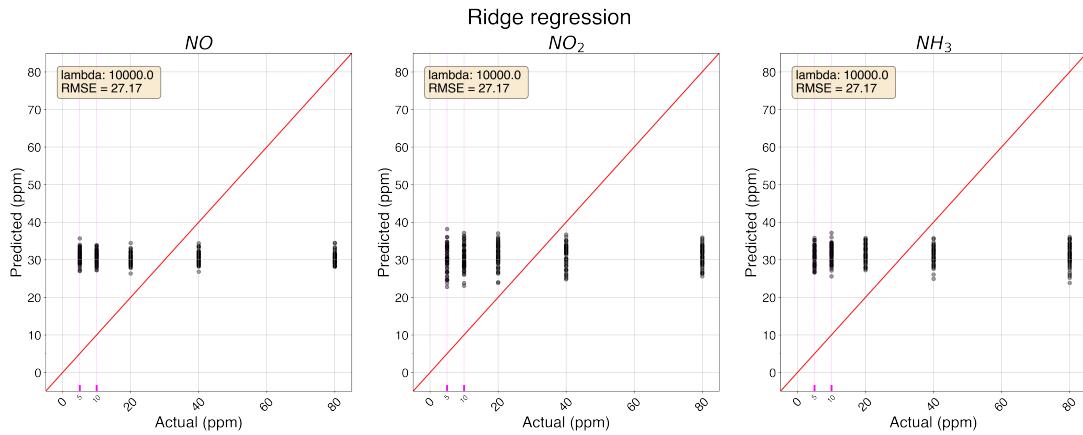


Figure 5.13: Ridge prediction for individual exposures

5.2.1 OLS - Average only

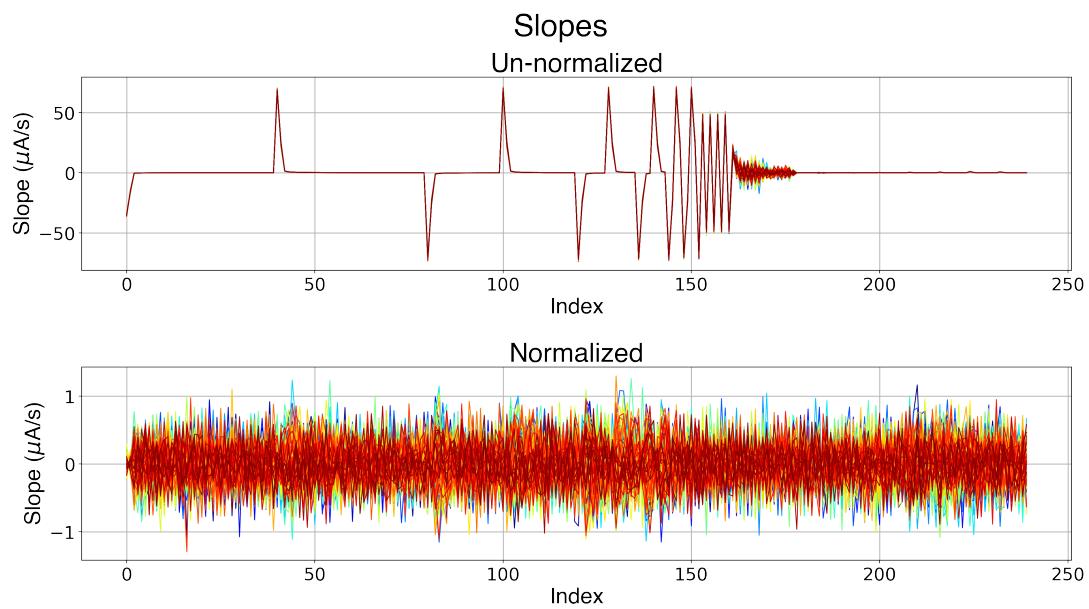


Figure 5.14: Slope features. Each line represents one of the 125 unique exposures.

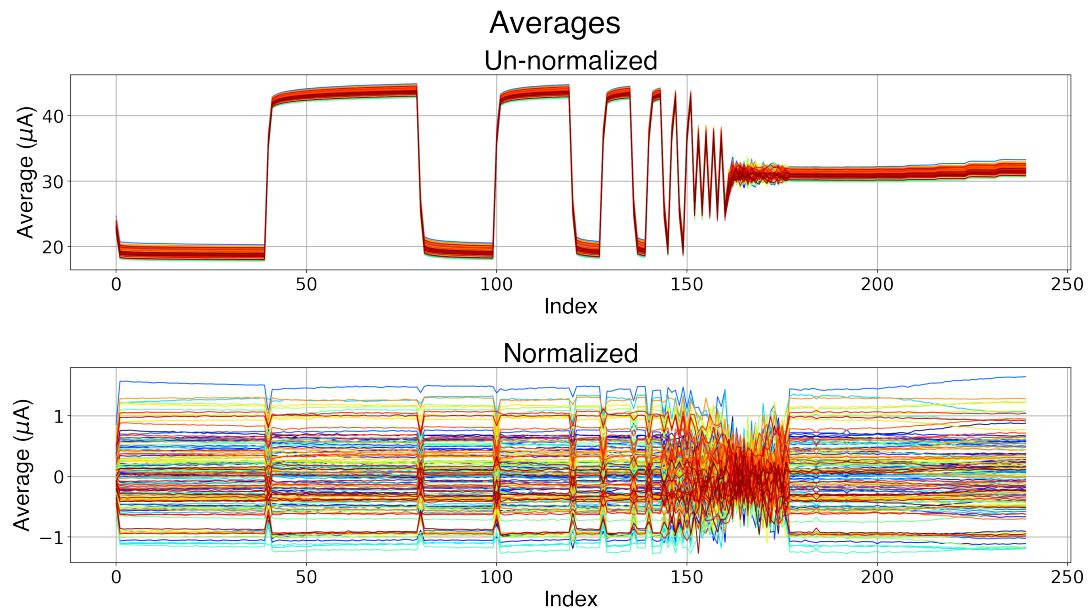


Figure 5.15: Average features. Each line represents one of the 125 unique exposures.

5.2.2 PCR - Average only

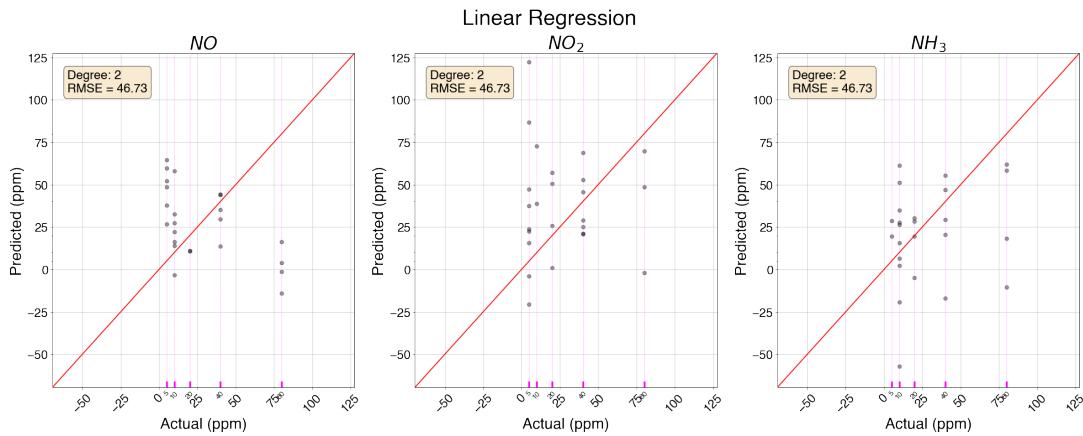


Figure 5.16: OLS prediction for unique gas mixtures.

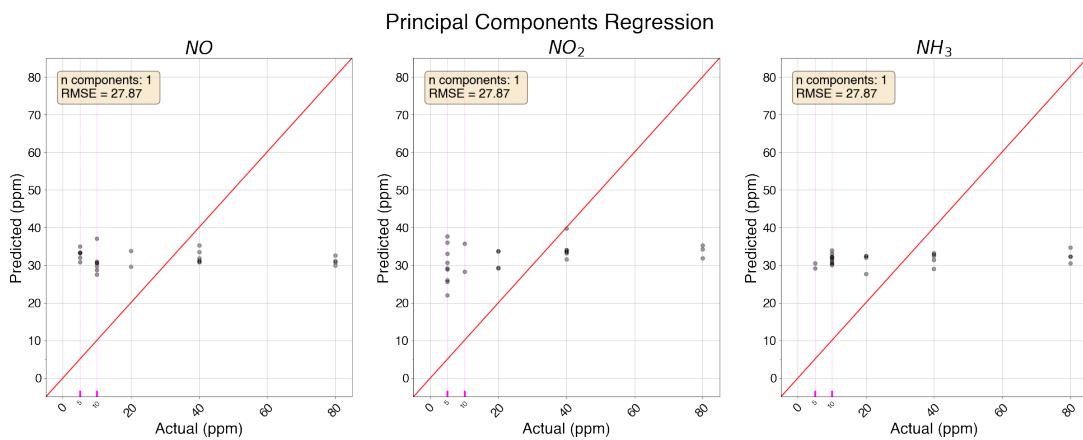


Figure 5.17: PCR prediction for unique gas mixtures.

5.2.3 PLSR - Average only

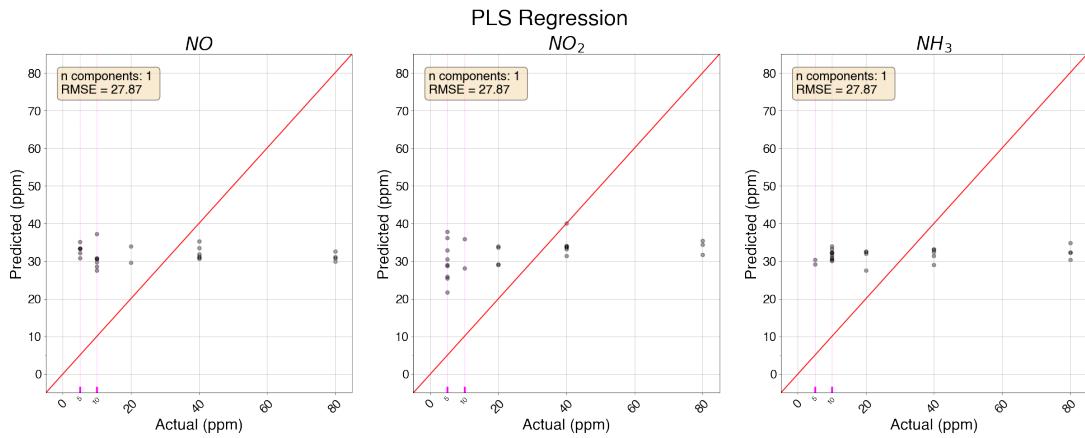


Figure 5.18: PLSR prediction for unique gas mixtures.

5.2.4 Ridge Regression - Average only

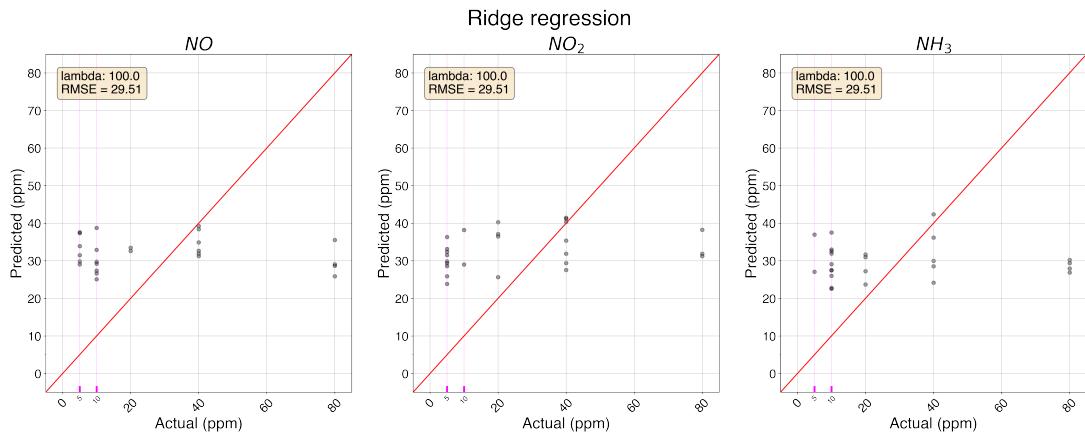


Figure 5.19: Ridge prediction for unique gas mixtures.



6 Discussion

This chapter is dedicated to explain the results obtained in the previous sections and relating them to statistical theory. The main objective here is to explain why the results are not satisfactory for gas concentration predictions.

6.1 Results

From the results shown in Chapter 5, it is clear that all models fail in predicting gas concentrations. As a first visual assessment, it can be seen from Figures 5.14 and 5.15 that there seems to be no clear order of response variables, i.e. simultaneous gas concentrations.

Slope features, for example, are approximately zero throughout the cycle, with the exception of some particular measurements as seen in Figure 5.14. Even then, visual inspection indicates that this feature is not informative of gas concentrations. For this reason, the second part of the analysis was done without these features.

Average features, on the other hand, seem to have some separation and indicate that gas concentrations might be explained by it. Nonetheless, it is not possible to order in any clear way. For example, mixtures with high concentrations of NO have average features that vary wildly, and it does not seem to follow any particular linear order of ammonia or NO_x concentration levels. Further attempts at ordering the data can be found in Appendix B. (**Note to Annika: I'll add those attempts we discussed (4 colors, etc.) in the appendix to avoid too much clutter. What do you think?**)

Ordinary Least Squares was expected to perform poorly, due to the high correlation of average features evidenced in Figure 5.1. Indeed, the results fail to predict gas concentrations at a reasonable level as seen on Figure 5.2. Looking at the actual vs. predicted plot, it can be seen that the predictions are centered at the mean of concentrations (31 ppm) and have high variance, indicated by the wide range of prediction in all concentration levels.

PCA with two components confirmed previous suspicions: there is no clear separation of gas concentrations for any of the gasses. Although Figure 5.4 shows that 80% of the variance can be explained by approximately 80 components, cross-validation indicate that only one PC yields minimal error before it starts to over fit the training data. Predictions in Figure 5.6 are poor, but have significantly less variance around the concentration mean than OLS. This is expected from this method, as the extraction of PCs is tightly related to the explained variability of the predictors, selecting linear combinations ordered by "importance" to the result.

PLSR, a method that has been shown to work in this type of problem, also performed poorly. Once again, there is not clear separation of concentration levels as shown in Figure 5.7, and CV shows that only one component, again, yields minimal RMSE. Prediction results in Figure 5.10 is very similar to predictions from PCR: centered around the mean with lower variance than OLS.

The final proposed model, Ridge regression, also fails completely in prediction, but brings meaningful insights for analysis. The CV plot for the shrinkage factor λ shows a curious behavior: with low values of λ , regression seems to fit training data well, with a relatively low RMSE of approximately 23. However, this is not the case for validation set. From the plot, extremely high values of regularization yields lowest RMSE in the test set, around 27. From Equations 3.22, 3.23 and 3.24, it can be seen that for high λ , the regression converges to a model that only predicts the mean (in this case, 31 ppm). A virtually "infinite" regularization would achieve best prediction in this case.

6.2 Method

In hindsight, the results indicate that the selection of linear models was ill-advised. Although PLSR seems to work well for problems using TCO, this is not the case for frequency modulation with NO_x and ammonia. For future work, non-parametric models are recommended. Perhaps their high flexibility and lack of assumptions about data could be of aid in achieving better prediction metrics.

Additionally, the frequency cycle itself could be changed. Most notably, instead of a square wave signal, a triangular wave could be more desired, as it would imply in less "stable" sections of the sensor response, possibly yielding more meaningful features, specially slopes.

6.3 The work in a wider context

From a statisticians point of view, it is always a misfortune when analysis' results are not "good" in the sense of high accuracy or low prediction error. These arbitrarily "bad" results, however, bring more insight to the problem, and knowing what does not work might be as valuable as knowing what works.

The quantification of NO_x and Ammonia, as shown in Chapter 1, is of paramount importance in the current world, where combustion processes are still commonplace. Although the advent of ever-improving electric vehicles is a silver lining regarding gas emissions and combustion processes, some industrial processes cannot avoid it.

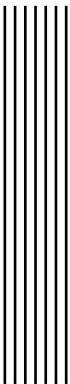


7 Conclusion

Given previous discussions and analysis, the answer to the research question "Can frequency modulation be used to simultaneously quantify NO_x and Ammonia concentrations?" seems to be: "perhaps not". Although poor predictions, the methods used here are far from exhausting the several other, possibly more flexible, models in the statisticians toolbox. Moreover, experiments using different frequency modulations (e.g. triangular waves instead of square) and/or different sensor calibrations and/or different temperatures could be further investigated in order to answer this question conclusively.

As for the second question "Which method yields best predictions of gas concentrations?", the correct answer would be "none". Nonetheless, out of all proposed methods, PLSR has best relative performance in terms of RMSE.

The author finds comfort in perhaps pointing future work towards better methods and possibly better quantification of these gases in hopes of addressing the problem more efficiently than current practices. In this sense, this thesis work is considered successful.



Bibliography

- Alberto Bernabeo, R., K. Webster, and M. Onofri (n.d.). "Health and Environmental Impacts of Nox: An Ultra- Low Level of Nox (Oxides of Nitrogen) Achievable with A New Technology." In: *Global Journal of Engineering Sciences* 3 (), pp. 2–7. DOI: 10.33552/gjes.2019.02.000540.
- ASTDR (2004). "Sheet for ammonia published by the Agency for Toxic Substance and Disease Registry (ASTDR)." In: 2672, pp. 1–18. URL: <https://www.atsdr.cdc.gov/MHMI/mmg126.pdf> <https://www.atsdr.cdc.gov/mmg/mmg.asp?id=7&tid=2#bookmark02>.
- Bastuck, Manuel (Jan. 2019). "Improving the performance of gas sensor systems with advanced data evaluation, operation, and calibration methods." PhD thesis, p. 267.
- Bishop, Christopher M (2006). *Pattern recognition and machine learning*. Springer.
- Boningari, Thirupathi and Panagiotis G. Smirniotis (2016). "Impact of nitrogen oxides on the environment and human health: Mn-based materials for the NOx abatement." In: *Current Opinion in Chemical Engineering* 13.x, pp. 133–141. ISSN: 22113398. DOI: 10.1016/j.coche.2016.09.004. URL: <http://dx.doi.org/10.1016/j.coche.2016.09.004>.
- Bur, Christian, Manuel Bastuck, Anita Lloyd Spetz, Mike Andersson, and Andreas Schütze (2014). "Selectivity enhancement of SiC-FET gas sensors by combining temperature and gate bias cycled operation using multivariate statistics." In: *Sensors and Actuators B: Chemical* 193, pp. 931–940. ISSN: 0925-4005. DOI: <https://doi.org/10.1016/j.snb.2013.12.030>. URL: <https://www.sciencedirect.com/science/article/pii/S0925400513015037>.
- Dunn, Kevin (2021). *Process Improvement Using Data*. McMaster University. ISBN: 9781292037578. URL: <https://learnche.org/pid/>.
- Forzatti, Pio (2001). "Present status and perspectives in de-NOx SCR catalysis." In: *Applied Catalysis A: General* 222.1. Celebration Issue, pp. 221–236. ISSN: 0926-860X. DOI: [https://doi.org/10.1016/S0926-860X\(01\)00832-8](https://doi.org/10.1016/S0926-860X(01)00832-8). URL: <https://www.sciencedirect.com/science/article/pii/S0926860X01008328>.

- Friedman, Jerome, Trevor Hastie, Robert Tibshirani, et al. (2001). *The elements of statistical learning*. Vol. 1. 10. Springer series in statistics New York.
- Guthrie, Susan, Sarah Giles, Fay Dunkerley, Hadeel Tabaqchali, Amelia Harshfield, Becky Ioppolo, and Catriona Manville (2018). *Impact of ammonia emissions from agriculture on biodiversity: An evidence synthesis*. Santa Monica, CA: RAND Corporation. DOI: 10.7249/RR2695.
- James, Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani (2013). *An introduction to statistical learning*. Vol. 112. Springer.
- Johnson, R.A. and D.W. Wichern (2013). *Applied Multivariate Statistical Analysis: Pearson New International Edition*. Pearson Education Limited. ISBN: 9781292037578. URL: <https://books.google.se/books?id=xCipBwAAQBAJ>.
- Ng, Kee Siong (2013). "A simple explanation of partial least squares." In: *The Australian National University, Canberra*.
- USEPA (2019). *Nitrogen Oxides Control Regulations*. <https://www3.epa.gov/region1/airquality/nox.html>. Accessed 2021-02-09.
- Wold, Svante, Michael Sjöström, and Lennart Eriksson (2001). "PLS-regression: a basic tool of chemometrics." In: *Chemometrics and Intelligent Laboratory Systems* 58.2. PLS Methods, pp. 109–130. ISSN: 0169-7439. DOI: [https://doi.org/10.1016/S0169-7439\(01\)00155-1](https://doi.org/10.1016/S0169-7439(01)00155-1). URL: <https://www.sciencedirect.com/science/article/pii/S0169743901001551>.
- Wright, Kevin (2017). *The NIPALS algorithm*. https://cran.r-project.org/web/packages/nipals/vignettes/nipals_algorithm.html. Accessed: 2021-03-12.



A

Data acquisition time stamps

Table A.1: Data acquisition timestamps.

Frequency (Hz)	Duration (s)	Feature	Start time (s)	End time (s)	Frequency (Hz)	Duration (s)	Feature	Start time (s)	End time (s)
0.05	20	Slope	0,0	0,4	25.0	2	Slope	44,0	44,4
		Average	9,6	10,0			Average	44,6	45,0
		Slope	10,0	10,4			Slope	45,0	45,4
		Average	19,6	20,0			Average	45,6	46,0
0.1	10	Slope	20,0	20,4	50.0	2	Slope	46,0	46,4
		Average	24,6	25,0			Average	46,6	47,0
		Slope	25,0	25,4			Slope	47,0	47,4
		Average	29,6	30,0			Average	47,6	48,0
0.25	4	Slope	30,0	30,4	100.0	2	Slope	48,0	48,4
		Average	31,6	32,0			Average	48,6	49,0
		Slope	32,0	32,4			Slope	49,0	49,4
		Average	33,6	34,0			Average	49,6	50,0
0.5	2	Slope	34,0	34,4	200.0	2	Slope	50,0	50,4
		Average	34,6	35,0			Average	50,6	51,0
		Slope	35,0	35,4			Slope	51,0	51,4
		Average	35,6	36,0			Average	51,6	52,0
1.0	2	Slope	36,0	36,4	500.0	2	Slope	52,0	52,4
		Average	36,6	37,0			Average	52,6	53,0
		Slope	37,0	37,4			Slope	53,0	53,4
		Average	37,6	38,0			Average	53,6	54,0
2.0	2	Slope	38,0	38,4	1000.0	2	Slope	54,0	54,4
		Average	38,6	39,0			Average	54,6	55,0
		Slope	39,0	39,4			Slope	55,0	55,4
		Average	39,6	40,0			Average	55,6	56,0
5.0	2	Slope	40,0	40,4	2500.0	2	Slope	56,0	56,4
		Average	40,6	41,0			Average	56,6	57,0
		Slope	41,0	41,4			Slope	57,0	57,4
		Average	41,6	42,0			Average	57,6	58,0
10.0	2	Slope	42,0	42,4	5000.0	2	Slope	58,0	58,4
		Average	42,6	43,0			Average	58,6	59,0
		Slope	43,0	43,4			Slope	59,0	59,4
		Average	43,6	44,0			Average	59,6	60,0