

# Prediction of ignition delay using data-driven framework for straight chain alkanes

Pragneshkumar Rana, Sivaram Ambikasaran, Krithika Narayanaswamy  
(Not sure about order: last considered as main guide and principal supervisor)

*Indian Institute of Technology, Madras*

---

## Abstract

Ignition delay is important global combustion property. Ignition delay time(IDT) is generally measured using Shock-tube and RCM experiments. Calculation of IDT from numerical reacting simulation is computationally costly and time consuming process. To obtain IDT faster and accurately, shock tube experimental data is used to predict IDT using machine learning algorithms. Generally, IDT correlation is Arrhenius type in its nature which reformulated using bond energy. Rather than activation energy, fuel bond information is used to avoid uncertainty and dependency of experimental parameters. To predict IDT, by k-means algorithm, dataset is divided into sub-dataset by Euclidean norm minimization. Optimum number of sub dataset(data clusters) is obtained using silhouette criteria. For each sub-dataset, using multiple regression and hypothesis testing final correlation and framework obtained to predict IDT for unknown fuel. Result obtained using framework and correlation shows excellent agreement with experimental result.

**Keywords:** Ignition delay prediction, Machine learning, Data-Driven , K-Means, Silhouette Criteria, IDT correlation, Framework

---

## 1. Introduction:

Combustion process is mainly characterized by transport processes and chemical reactions. When fluid undergoes chemical reaction, it liberates heat without external source of energy such that process sustains, which is called as Ignition. Ignition comprises series of coincidental physical and chemical processes which have different characteristic time scale, which is called as ignition delay.

Ignition delay is physio-chemical property. It is one of the major global combustion property. Ignition delay gives important information about fuel reactivity and ignition. Ignition delay is mainly comprised of two parts: physical ignition delay and chemical ignition delay. Physical ignition delay depends on certain physical phenomena such as heating, fuel atomization, penetration of spray, and evaporation rate of fuel for different temperature range. Whereas, chemical ignition delay is mainly function of chemical characteristics of fuel, molecular structure, equivalence ratio, etc. The main focus of this study is chemical ignition delay. Henceforth, chemical ignition delay will be referred as Ignition delay(ID). Apart from chemical characteristics of fuel, Ignition delay is also function of pressure, temperature as reaction rates are also function of same.

Ignition delay is crucial factor for design of combustor. Right amount of ignition delay time is required for proper combustion. Prediction of IDT for various fuel and wide range of conditions is

complicated process. Ignition delay are generally calculated using reacting flow simulation which involves large number species and thousands of reaction including broad range of chemical and flow time scale. Prediction of IDT using realistic fuel chemistry requires large scale simulations of combustion phenomena [1]. Ignition delay solution using realistic-detail mechanism requires full-scale chemistry solvers which are computationally intensive. Taking motivation from such complications, acquired framework and correlations are simplified, accurate, and efficient which is applicable over variety of fuels and wide range of conditions.

For development and validation of predictive reaction models, ignition data are very useful. Substantial work has been done by researchers to calculate and correlate the ignition delay. Major ignition delay equations are based on Arrhenius-type correlation which can be written as [2]:

$$\tau = A\phi^\alpha P^\beta X_{O_2}^\gamma \exp(\lambda) \quad (1)$$

$\alpha, \beta, \gamma$  are polynomials, to capture changes in functionality across different regimes.  $\lambda$  is overall activation energy. Horning et al. [3] has conducted study of different hydrocarbon at high-temperature to observe auto ignition and thermal decomposition. He found that for n-alkanes at  $\phi = 1$  follows the given ignition delay correlation-2

$$\tau = 9.40 * 10^{-6} P^{-0.55} X_{Oxi}^{0.63} C^{-0.5} e^{46500/RT} \quad (2)$$

Where C is number of carbon atoms in the molecules. Such Arrhenius-type equations are constrained based on certain physical condition. In parametric uncertainty [3], it was mentioned that activation energy is very sensitive to ignition temperature which directed to find out more efficient parameter to replace activation energy.

#### Notations :

P = Pressure	T = Temperature
$\phi$ = Equivalence Ratio	$\tau$ = Ignition Delay Time
$X_{Oxi}$ = Oxidizer mole fraction	$X_{Fuel}$ = Fuel mole fraction
$E_a$ = Activation Energy	$\delta H$ = Change in enthalpy
R = Gas constant	A = Scaling factor
IDT = Ignition delay time	C = Constant

In recent study of distillate fuels, Fethi and Amir [4] has obtained ignition delay correlation for gasoline and Jet-Fuel using modified Arrhenius expression which applicable over wide range of conditions [ $P = 10\text{-}80$  bar,  $P_0 = 1$  bar,  $\phi = 0.5\text{-}2$ , fuel/air mixtures, units are ms, bar, K, mol, kcal]. Using chemical model of Sarathy et al. and by montecarlo simulation they generated  $10^7$  samples.

Using those samples and by statistical approach they found the correlation [3]:

$$\begin{aligned}\tau_{gasoline} &= 6.76 * 10^{-7} \frac{P^{-1.01}}{20} \phi^{1.13 - \frac{(17.59)}{T}} \exp \frac{29.39}{RT} \\ &\quad \text{for } T > \frac{1000}{-0.073 \ln(\frac{P}{P_0}) + \phi^{-0.0338} + 0.0938} \\ \tau_{JetFuel} &= 4.76 * 10^{-7} \frac{P^{-1.21}}{20} \phi^{2.04 - \frac{(29.56)}{T}} \exp \frac{29.33}{RT} \\ &\quad \text{for } T > \frac{1000}{-0.0371 \ln(\frac{P}{P_0}) + \phi^{-0.00727} + 0.0995}\end{aligned}\tag{3}$$

Such correlations are useful for modelling of the surrogates. But such correlations are fuel specific and constrained by physical conditions so it is quite difficult to obtain IDT for new fuel or surrogates. **The goal** of present study is to find out correlation which generalize, efficient, applicable over wide range physical condition and variety of fuel. To obtain this, it is necessary to modify Arrhenius-type equation. In next section, ignition delay formulation has been discussed which is further used to find out coefficient using machine learning approach.

### 1.1. Ignition delay formulation:

As we discussed earlier, many researchers have obtained IDT correlation from Arrhenius-type formulation. But specificity of fuel and physical parameter causes uncertainty in result of IDT. To remove such complication, ignition delay correlation has to be reformulated.

In chemical reactions, the most sensitive parameter in ignition delay correlation is Activation energy. Activation energy describes overall transformation of reaction and it only gives macroscopic information about reaction as intermediates are not considered in any reaction. More microscopic information about the single step reaction can be obtained using the Eyring equation [5].

$$k = \kappa \left( \frac{k_B T}{h} \right) e^{\frac{\Delta S^\ddagger}{R}} e^{\frac{-\Delta H^\ddagger}{RT}}\tag{4}$$

Where,  $k_B$  = Boltzmann constant,  $T$  = absolute temperature,  $h$  = Planck's constant,  $\Delta S^\ddagger$  = Activation Entropy,  $\Delta H^\ddagger$  = Activation Enthalpy. Eyring's equation is based on statistical and mechanical rationale of transition state theory whereas Arrhenius equation is empirical. These both equation are different in its nature. Relation between these two equations is possible when elementary reaction is as uni-molecular or bi-molecular. In such case, activation energy or Energy barrier can be defined in terms of enthalpy of activation [5] which is closely related to bond energy.

$$E_a = \Delta H^\ddagger + nRT\tag{5}$$

According to transition state theory, when molecules with enough kinetic energy collides in certain orientation, it may generate activated complex. Bond structure of activated complex is different from reactants bond structure.  $\Delta H^\ddagger$  plays critical role in bond formation or breakage. So, enthalpy of activated complex is related to enthalpy of reaction.

In combustion, heat of combustion ( $\Delta H_{combustion}$ ) or heat of reaction ( $\Delta H_{reaction} = -\Delta H_{combustion}$ )

is directly related to bond dissociation energy [6] which can be expressed as,

$$\begin{aligned}\Delta H_{combustion} &= H_{reactants} - H_{products} \\ &\approx \text{Bond energy of Reactants} - \text{Bond energy of Products}\end{aligned}\quad (6)$$

Being point function, enthalpy of activation ( $\Delta H^\ddagger$ ) for forward reaction can be expressed as difference between reactant state and activation state,

$$\Delta H^\ddagger = H_{reactant} - H^\ddagger \quad (7)$$

from relation-(6),

$$\Delta H^\ddagger = \Delta H_{combustion} + H_{products} - H^\ddagger \quad (8)$$

In IDT relation-(2) Horning et al. has showed that IDT depends on the number of carbons. So, Arrhenius-type ignition delay formulation can be written as below using equation-(5) and (8).

$$\begin{aligned}\tau &= A \cdot \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{E_a}{RT}\right) \\ &= A \cdot \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\Delta H^\ddagger + nRT}{RT}\right) \\ &= A \cdot \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\Delta H_{combustion} + H_{products} - H^\ddagger}{RT}\right) \exp\left(\frac{nRT}{RT}\right) \\ &= A \cdot \exp(n) \cdot \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\Delta H_{combustion} + H_{products} - H^\ddagger}{RT}\right) \\ &= C \cdot \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\Delta H_{combustion} + H_{products} - H^\ddagger}{RT}\right)\end{aligned}\quad (9)$$

From-6, it clear that enthalpy of combustion depends on the type of bond and bond energy. Other enthalpy details are constant being point function so, formulation can be rewritten as,

$$\begin{aligned}&\propto \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\text{Fuel Bond Energy} \cdot \text{Constant}}{T}\right) \\ &\propto \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\sum_{i=1}^n (\text{Bond Type})_i \cdot (\text{Bond Energy})_i \cdot \text{Constant}}{T}\right) \\ &\propto \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\sum_{i=1}^n (\text{Bond Type})_i \cdot (\text{Bond Energy} \cdot \text{Constant})_i}{T}\right) \\ &\propto \phi^\alpha \cdot P^\beta \cdot X_{O_2}^\gamma \cdot \exp\left(\frac{\sum_{i=1}^n (\text{Bond Type})_i \cdot (\text{Bond Constant})_i}{T}\right)\end{aligned}\quad (10)$$

44 Bond constant has unit of K/J which is inverse of universal gas constant. This formulation gives  
45 indication that ignition delay depends on the chemical bonds of the fuel. Further in the formulation,  
46 to make quantities independent of unit, it was normalized with standard conditions. Considering  
47 all the affecting parameters in above formulation-10 hypothetical functional form of IDT can be

48 written as,

$$\tau = f(T, P, \phi, X_{Fuel}, X_{O_2}, X_{Dilutant}, \text{Type of bonds}) \quad (11)$$

49 But it is known that ,

$$X_{Fuel} + X_{O_2} + X_{Dilutant} = 1 \quad \& \quad \phi = \frac{\left(\frac{X_{fuel}}{X_{O_2}}\right)_{act}}{\left(\frac{X_{fuel}}{X_{O_2}}\right)_{stochio}} \quad (12)$$

As Fuel composition is already known, from stoichiometry equation  $\left(\frac{X_{fuel}}{X_{O_2}}\right)_{stochio}$  is attainable. It is possible to obtain  $\phi$  using  $X_{fuel}$  and  $X_{O_2}$ . Same is true for  $X_{Dilutant}$ . Removing redundant parameters formulation can be written as,

$$\tau = f(T, P, X_{Fuel}, X_{O_2}, \text{Type of bonds}) \quad (13)$$

$$\begin{aligned} \tau &\propto \left(\frac{T}{T_0}\right)^a \cdot \left(\frac{P}{P_0}\right)^b \cdot X_{Fuel}^c \cdot X_{O_2}^d \cdot \exp\left(\frac{\sum_{i=1}^n (\text{Bond Type})_i \cdot (\text{Bond coefficient})_i}{\frac{T}{T_0}}\right) \\ &= C' \cdot \left(\frac{T}{T_0}\right)^a \cdot \left(\frac{P}{P_0}\right)^b \cdot X_{Fuel}^c \cdot X_{O_2}^d \cdot \exp\left(\frac{\sum_{i=1}^n (\text{Bond Type})_i \cdot (\text{Bond coefficient})_i}{\frac{T}{T_0}}\right) \end{aligned} \quad (14)$$

50 formulation can simplified by taking natural log on both side,

$$\begin{aligned} \ln(\tau) &= \ln(C') + a \cdot \ln\left(\frac{T}{T_0}\right) + b \cdot \ln\left(\frac{P}{P_0}\right) + c \cdot \ln(X_{Fuel}) \\ &\quad + d \cdot \ln(X_{O_2}) + \left(\frac{\sum_{i=1}^n (\text{Bond Type})_i \cdot (\text{Bond coefficient})_i}{\frac{T}{T_0}}\right) \end{aligned} \quad (15)$$

$$\begin{aligned} \ln(\tau) &= C + a \cdot \ln\left(\frac{T}{T_0}\right) + b \cdot \ln\left(\frac{P}{P_0}\right) + c \cdot \ln(X_{Fuel}) \\ &\quad + d \cdot \ln(X_{O_2}) + \sum_{i=1}^n \left( (\text{Bond coefficient})_i \cdot \frac{T_0 \cdot (\text{Bond Type})_i}{T} \right) \end{aligned} \quad (16)$$

51 In this equation, C,a,b,c,d are coefficients which can be obtained from ignition delay data using  
52 multiple regression. It was assumed that ignition delay correlation depends on all the parameters  
53 which was refined using hypothesis testing. To obtain correlation, multiple regression along with  
54 machine learning algorithm was used. Zeroth step of whole process was to collect the data. Data is  
55 a heart for machine learning technique. Further discussion is about data collection and processing.

Fuel		Temperature (K)	Temperature Uncertainty (%)	Pressure (atm)	Pressure Uncertainty (%)	Fuel Mole Fraction	Oxygen Mole Fraction	Equivalence Ratio	Data Points	Research Group	Reference
Ethane	max	2497	$\pm 2.3$	20	$\pm 0.7$	2	7	2	134	Sanford, Xi'an Jiaotong	[7],[8],[9],[10]
	min	1086	$\pm 0.7$	0.537	$\pm 0.7$	0.01	0.0972	0.35			
Propane	max	1841	$\pm 3$	67.8	$\pm 1$	4	20	5	174	Stanford, Xi'an Jiaotong	[11],[12],[13],[14],[15],[16],[9]
	min	950	$\pm 0.7$	1.12	$\pm 0.7$	0.05	0.25	0.5			
Butane	max	1761	$\pm 3$	5.5	$\pm 1$	2	13	2	58	Stanford	[12],[11],[13],[9]
	min	1230	$\pm 0.7$	1.03	$\pm 0.7$	0.05	0.325	0.5			
Pentane	max	1533	$\pm 0.7$	3.75	$\pm 1$	0.5	4	1	15	Stanford	[14],[17]
	min	1261	$\pm 0.7$	1.62	$\pm 1$	0.25	4	0.5			
Hexane	max	1475	$\pm 0.7$	3.6	$\pm 1$	0.42	4	1	16	Stanford	[14],[17]
	min	1237	$\pm 0.7$	1.67	$\pm 1$	0.21	4	0.5			
heptane	max	1784	$\pm 1.8$	60.6	$\pm 1$	1.874	20.6	2	107	Stanford	[11],[12],[13],[18],[19],[20],[21]
	min	806	$\pm 0.7$	1.14	$\pm 0.7$	0.03	0.33	0.5			
Octane	max	1455	$\pm 0.7$	3.81	$\pm 1$	0.32	4	1	15	Stanford	[14],[17]
	min	1252	$\pm 0.7$	1.87	$\pm 1$	0.16	4	0.5			
Nonane	max	1301	$\pm 0.7$	41.76	$\pm 1$	0.4	4	2	27	Stanford	[14],[17]
	min	1051	$\pm 0.7$	13.52	$\pm 1$	0.2	4	0.5			
Decane	max	1706	$\pm 2$	5.15	$\pm 0.7$	2.567	21	1.89	25	Stanford	[11],[12],[13],[22]
	min	1081	$\pm 1.5$	1.22	$\pm 9.6$	0.03	0.3875	0.64			
Dodecane	max	1657	$\pm 1$	33.7	$\pm 9.3$	2.138	21	1.88	162	Stanford	[23],[24],[25],[26],[27],[28],[29]
	min	727	$\pm 0.7$	2.07	$\pm 0.7$	0.0371	0.731	0.05			
Hexadecane	max	1355	$\pm 2$	6.77	$\pm 0.7$	0.1832	4	1.22	18	Stanford	[30],[31],[22]
	min	1159	$\pm 2$	1.71	$\pm 0.7$	0.0312	1	0.56			

Table 1: Summary of obtained fuel data from different research group along with physical condition and uncertainty in the physical parameter

## 2. Data collection and processing

Many researchers have performed experiment to measure ignition delay using different fuels for wide range of conditions. Focus of this study is to predict IDT of straight chain alkanes. Source of data points are literature and publications. Main source of Alkanes data is Stanford group and Xi'an Jiaotong group. Summary of Ignition delay data obtained for wide range of conditions with variety of fuels is given in table-1:

Xi'an Jiaotong group has performed Ethane and Propane experiments using shock-tube of diameter 11.5cm. They have collected data by measuring OH species profile. Due to smaller diameter, uncertainty in result is slightly higher.

Stanford group has performed IDT experiment for short-chain to long-chain alkanes using different shock-tubes. In experiments, IDT is mainly obtained using species profile of OH, CO<sub>2</sub>, CH, and CH<sub>3</sub> which is also a more accurate way to measure IDT. While measuring IDT from shock-tubes, inherent uncertainty such as boundary layer growth, reduction in amplitude of shocks takes place which is inevitable as such phenomena generates uncertainty in measurement. Uncertainty is also reported in the Table-1 which also useful for generation of data points. Gauthier et al. reported that ignition delay time obtained by experiments have average uncertainty of  $\pm 15\%$  [18]. In other literature, reported uncertainty value varies from  $\pm 10\%$  to  $\pm 30\%$ .

For certain data points, temperature and pressure uncertainty value was not reported which was replaced by the least available uncertainty value from whole dataset. It is clear from the dataset that, extensive study of Propane, Heptane, Dodecane has made availability lot of data points whereas Pentane, Hexane, Nonane and Decane has less experimental data points. Multi-linear regression from such imbalance dataset will cause bias in the result. Such issues were resolved by generating more data points using uncertainty value of each parameter. Procedure is explained in further discussion.

Apart from information about all physical parameter, fuel bond information is also useful for

IDT correlation. In next section, procedure to obtain fuel bond information is discussed.

### 3. Fuel Bond Information

Type of bonds in fuel plays key role in combustion process as chemical transformation is possible through bond breakage. Type and number of bonds varies with fuel.

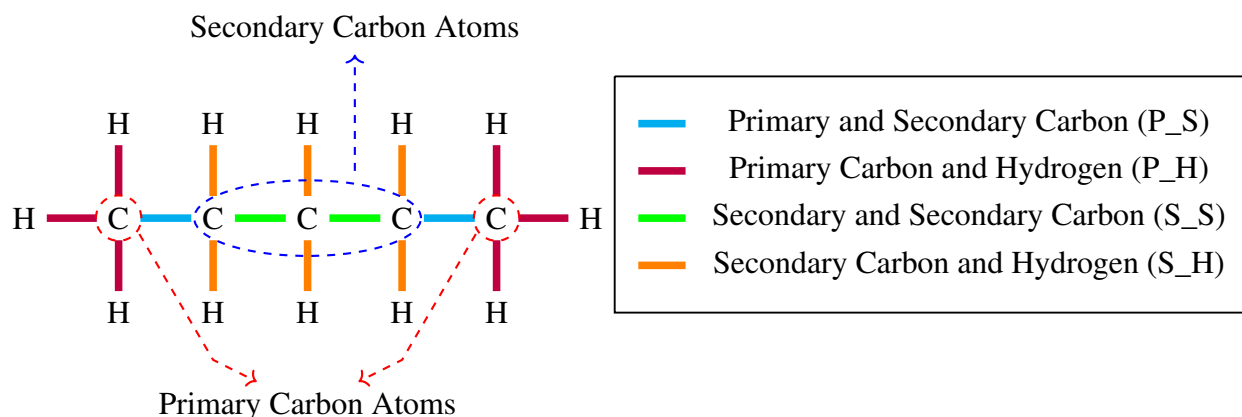


Figure 1: Different type of bonds in straight alkane (Pentane)

Bond information of fuel is obtained using SMILE (Simplified molecular-input line-entry system), which is method for computer to extract information and make it processable. It is based on molecular graph theory. SMILES follows certain rules to make entry of chemicals [32]. Chemical molecular structure are well defined by line entry system to acquire required information.

RDKit [33] in python was used to extract the bond information. Using SMILES and RDKit, Mol-file was obtained. Mol-file gives information about Bond connections and orientation of bond in 3D space. By text processing of Mol-Block(file) bond information was extracted.

It is clear from the figure-1 that, straight chain alkane (for carbon chain  $> 2$ ) contains 4 different type of bonds. Ethane is only fuel in which primary-primary carbon bond exist so that fuel is excluded from the analysis. Bond details were extracted from the Mol-file. Different type of bond detail are used as parameters in the dataset. Using all the obtain information final dataset given in table-2.

### 4. Data Clustering and Analysis:

The obtained ignition delay has wide range of conditions for temperature, pressure, Fuel and Oxygen Mole fraction for variety of fuel. Multi-linear regression of dataset without any treatment gives poor coefficient of determination  $R^2 = 0.79$  which indicates that dataset has to divide based on certain pattern.

Ignition delay data of Alkanes have 8 parameters. Visualization of such data points in higher dimension is not possible whereas straightforward in less dimensions. Weiqi et al. has divided the dataset based on the turnover states to reveals distinct thermodynamic and kinetic property under different pressure [34]. To develop IDT correlation for Hydrogen/air mixture, Zhao et al. [35] has

Fuel	T(K)	T_Error(%)	P(atm)	P_Error(%)	Fuel(%)	Oxygen(%)	P_S	S_S	P_H	S_H	Time( $\mu$ s)
CCC	1376	$\pm 1$	1.19	$\pm 1$	4	20	2	0	6	2	357
CCCC	1409	$\pm 1$	1.17	$\pm 1$	1	6.5	2	1	6	4	390
CCCCC	1395	$\pm 0.7$	3.47	$\pm 1$	0.5	4	2	2	6	6	316
CCCCCC	1273	$\pm 0.7$	3.32	$\pm 1$	0.42	4	2	3	6	8	1046
CCCCCCC	1378	$\pm 1$	2.326	$\pm 1$	0.03	0.33	2	4	6	10	1330
CCCCCCCC	1289	$\pm 0.7$	2.01	$\pm 1$	0.32	4	2	5	6	12	1198
CCCCCCCCC	1107	$\pm 0.7$	14.8	$\pm 1$	0.4	4	2	6	6	14	965
CCCCCCCCC	1081	$\pm 2$	5.14	$\pm 9.6$	1.44	21	2	7	6	16	1696
CCCCCCCCC	1045	$\pm 1$	6.71	$\pm 9.3$	0.6098	21	2	9	6	20	2753
CCCCCCCCC	1181	$\pm 2$	2.13	$\pm 1$	0.1776	4	2	13	6	28	5536

Table 2: Fuel dataset made after obtaining bond information. Uncertainty in parameter is replace by least available value in-case not reported. Fuel is represented using SMILE. T\_Error(%) and P\_Error(%) represents uncertainty in measurement of temperature and pressure respectively. Bond notation are same as mentioned in figure-1

divided the dataset into six sub domains based on high, intermediate and low pressure along with low and High temperature range. Such finding motivates to divide the dataset into sub parts.

Due to recent advancement in machine learning, it is possible to unveil concealed pattern of n-Dimensional data using different clustering technique, which is one of the famous unsupervised learning technique for unlabelled data. Out of many available algorithm, K-Means is famous and preferred partition based algorithm. It performs better in terms accuracy and execution time compared to other algorithm[36]. K-Means algorithm generates K centroids and minimizes the distance between cluster centroids and data-points for given number of clusters.

#### K-Means Algorithm:

- Consider the data points  $x_i = x_1, x_2, x_3, \dots, x_n \in \mathbf{R}^d$  has j features.  $\mu_i \in \mathbf{R}^d$  denotes number of clusters and  $q_1, \dots, q_n \in \{1, \dots, K\}$  denotes points assigned to the centroid. Such that the sum of distances is minimized within cluster which is defined as,

$$E(\mu_1, \mu_2, \dots, \mu_K, q_1, q_2, \dots, q_n) = \sum_{i=1}^n ||x_i - \mu_{q_i}||_p^p \quad (17)$$

where p denotes the norm.

- steps:
  1. Randomly select the K initial cluster centroid.
  2. Assign the data points to nearest cluster
  3. Recompute the cluster based current data members of associated cluster
  4. Repeat the procedure till convergence criteria not met or clusters are not moving (Go to step:2)

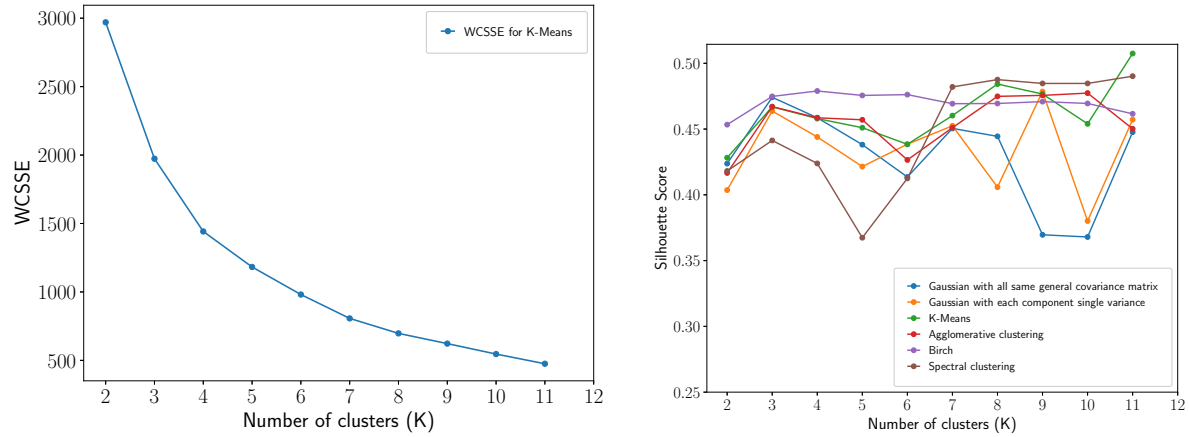
K-Means algorithm is used from scikit-learn library. Which computes the clusters using Elkan's algorithm [37]. K-Means is useful method when number of clusters are known. But for the case of IDT, number of clusters are unknowns. To find out optimal number of clusters, elbow method is used along with silhouette criteria.



Elbow method calculates within-cluster sum of squared error (WCSSE) for given k cluster by equation-18. The K-value for which WCSSE starts declining, gives the optimized number of clusters as it minimizes the WCSSE.

$$WCSSE = \sum_{i=1}^n \sum_{x \in S_i} ||x_i - \mu_i||^2 \quad (18)$$

For ignition delay data, WCSSE is calculated for 2 to 11 number of clusters. Apart from ethane, all available data has been used to generate the clusters. From fig-2(a), it was observed that elbow(sharp slope) was generated around K=3 or K=4. The obtained result from fig-18 is quite ambiguous to analyse. To have more robust decision Silhouette criteria was needed along with elbow method.



(a) Elbow Method Result - Within-cluster sum of squared error(WCSSE) vs Number of cluster to decide the optimum number of cluster for K-Means

(b) Silhouette Score vs Number of cluster for different clustering technique

Figure 2: Criteria to decide the number of clusters for all fuel components

### Silhouette Score Calculation:[38] [39]

- Data were divided into K-clusters
- Each data point  $x_i$  is assigned to cluster  $C_k$  such that  $\forall x_i \in C_k$
- Silhouette criteria measures the relative separation distance between points in same cluster to other clusters. Mean distance **within cluster** from point i to other points were measured by,

$$a(i) = \frac{1}{|C_i|-1} \sum_{j \in C_i, i \neq j} d(i, j) \quad (19)$$

where,  $d(i, j)$  is the distance between data points i and j in the cluster  $C_i$ . Numerator contains  $|C_i|-1$  as  $d(i, i)$  is excluded. In summary,  $a(i)$  shows how well data point i is assigned to its cluster. Smaller value is expected for proper assignment.

- Mean distance between point  $i$  in cluster  $C$  to **points in other cluster** were measured by following equation:

$$b(i) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j) \quad (20)$$

Use of min operator on cluster  $k$  shows that out of all value of different cluster, minimum value is used for calculation of  $b(i)$  which is nearest neighbouring cluster to that specific point- $i$ . Large value of  $b(i)$  suggest that point- $i$  does not properly matches with neighbouring cluster.

- Silhouette for points- $i$  is defined as,

$$s(i) = \begin{cases} \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} & |C_i| > 1 \\ 0 & |C_i| = 1 \end{cases} \quad (21)$$

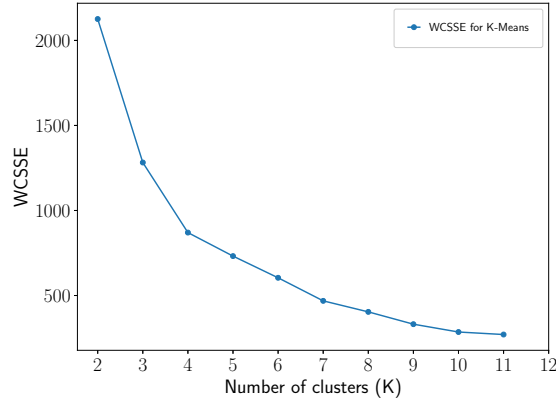
which is also defined as,

$$s(i) = \begin{cases} 1 - \frac{a(i)}{b(i)} & \text{if } |a_i| < |b_i| \\ 0 & \text{if } |a_i| = |b_i| \\ \frac{b(i)}{a(i)} - 1 & \text{if } |a_i| > |b_i| \end{cases} \quad (22)$$

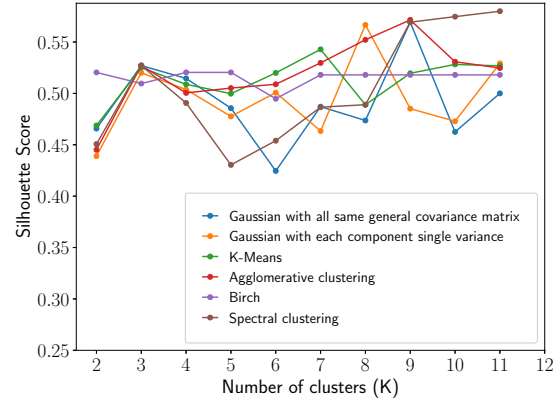
- $s(i)$  was calculated for all the points. The average value  $s(i)$  in the cluster shows how tightly all point are assembled in the cluster. Mean  $s(i)$  of of whole dataset shows how properly data points forms the group.

Silhouette analysis is useful for each data point to measure the distance between the clusters they belong and other neighbouring clusters. Silhouette Score varies from -1 to 1 in which 1 is ideal which shows that object are well separated and -1 shows improper clustering. To calculate Silhouette Score, `sklearn.metrics.silhouette_score` function is utilized.

For IDT dataset, Silhouette analysis has been done using different clustering algorithm. First drop in result was observed at cluster  $K=3$  for clustering techniques which emphasis that data should be divided into three partition. IDT dataset-1 has large variation in data points. Observed result is quite sensitive to type of fuel, conditions affecting to IDT and number of data points. For verification of technique and to check biasness in the result, silhouette analysis has been done on fuel which has more than 100 data points. From table-1, it observed that propane, heptane and dodecane satisfies that criteria. Those fuel also covers wide of range of conditions. The observed result given in fig-3(a) and fig-3(b) shows similar trend obtained in case if 'all fuel' - 2(b) 2(a). It clear from the fig-3(b) that number of cluster required for dataset of propane, heptane and dodecane is also 3.

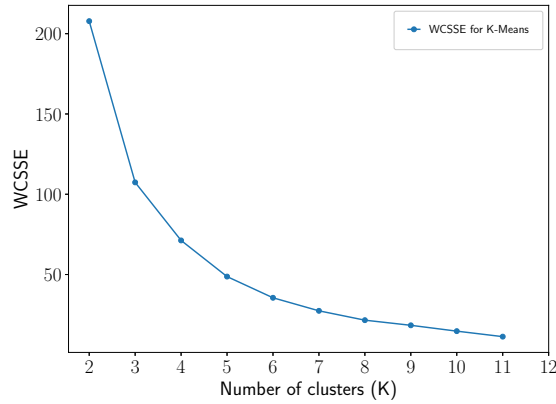


(a) Elbow Method Result: Within-cluster sum of squared error(WCSSE) vs Number of cluster for K-Means using data points of Propane, Heptane and Dodecane Fuel

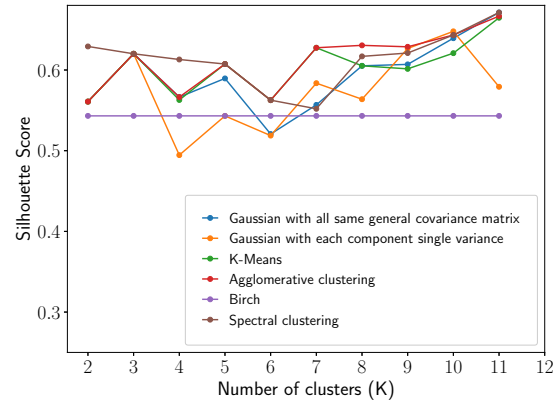


(b) Result of Silhouette Score vs Number of cluster for different clustering technique using data points of Propane, Heptane and Dodecane

Figure 3: Criteria to decide the number of clusters obtained using Propane, Heptane and Dodecane



(a) Elbow Method Result: Within-cluster sum of squared error(WCSSE) vs Number of cluster for K-Means using data points of Propane

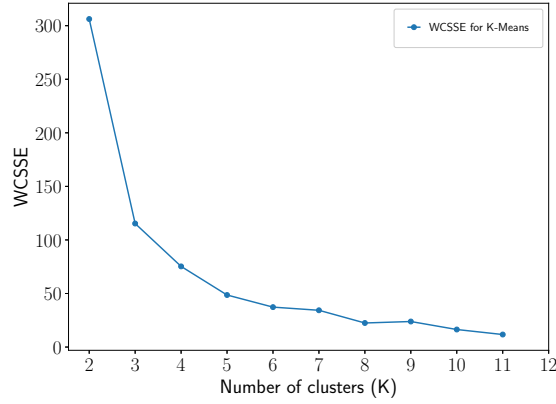


(b) Result of Silhouette Score vs Number of cluster for different clustering technique using data points of Propane

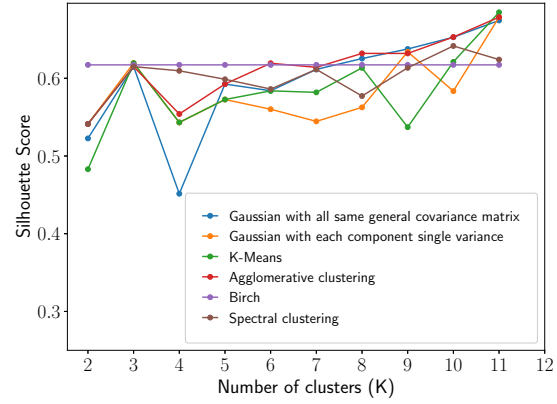
Figure 4: Criteria to decide the number of clusters obtained using Propane

Similar result was also observed in case of individual fuel dataset of Propane, Heptane as given in figure-4(b),5(b) which also agrees with expected trend. For propane and heptane observation matches with expected trend whereas in case of heptane -21(b), it suggest to divide the dataset into two sub-dataset as Dodecane data points has limited range of physical conditions. Dodecane data, in table-1 shows that, compared to heptane and propane data is limited by pressure at 33.7 whereas, propane and heptane has pressure around 60 atm. Similar observation is also made in equivalence ratio.

Silhouette score obtained in all the cases is above 0.4 which emphasis that data points are grouped by certain parameters. At a same time, clustering value is far away from 1 which means that they do not generate perfect cluster and data are quite separated which is useful information for regression analysis.

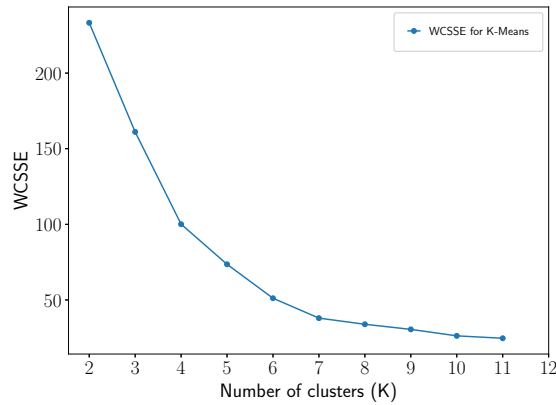


(a) Elbow Method Result: Within-cluster sum of squared error(WCSSE) vs Number of cluster for K-Means using data points of Heptane

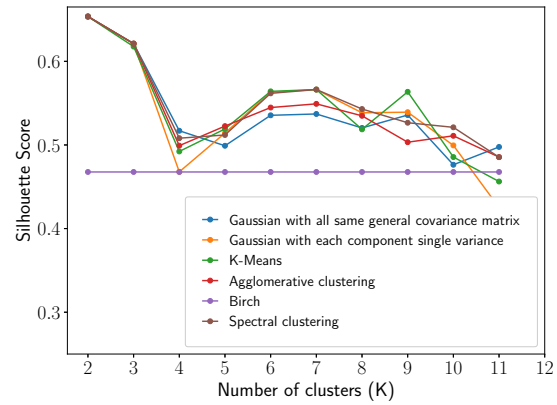


(b) Result of Silhouette Score vs Number of cluster for different clustering technique using data points of Heptane

Figure 5: Criteria to decide the number of clusters obtained using Heptane



(a) Elbow Method Result: Within-cluster sum of squared error(WCSSE) vs Number of cluster for K-Means using data points of Dodecane



(b) Result of Silhouette Score vs Number of cluster for different clustering technique using data points of Dodecane

Figure 6: Criteria to decide the number of clusters obtained using Dodecane

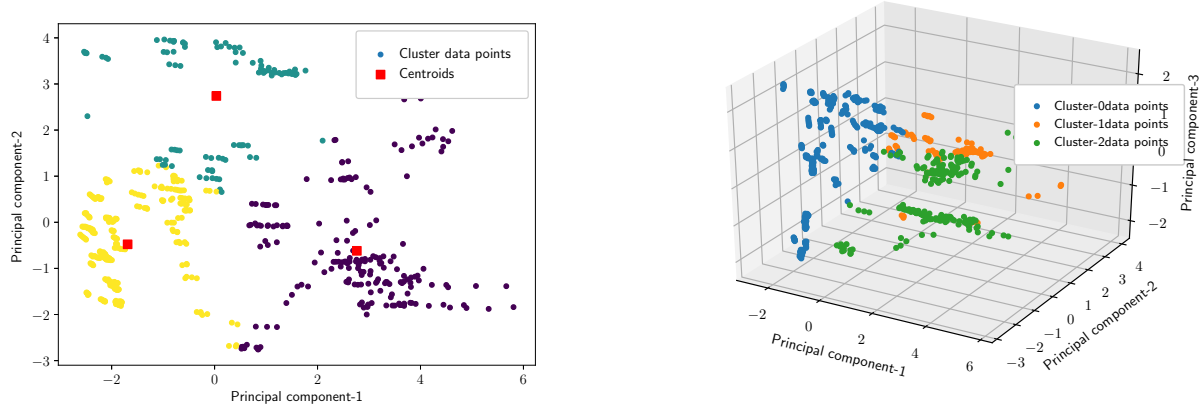
#### 4.1. Principal Component Analysis (PCA) of Dataset

Principal component analysis is useful for dimensionality reduction and for analysis of major component of the data. For given data, PCA find out major eigen values and eigen vector. Principle components obtained using full SVD solver showed that, data has 6 major eigen values and other two are near to 0. For (616, 8) size of dataset obtained variation in principal components(Eigen values) are obtained as below:

$$\begin{matrix} 4.80958553e+00 & 2.55919622e+00 & 1.15728172e+00 & 9.98961156e-02 \\ 2.33198384e-02 & 1.68372670e-04 & 2.93116344e-32 & 1.18120890e-33 \end{matrix}$$

From six eigenvalues it is clear that, it has major three principal eigen values and out of rest, other three are smaller eigen values. so, it is possible to visualize the major variation in data and

clustering in the 3-Dimensions. To visualize the clustering of data in 3D, K-Means along with PCA was implemented. Obtained result are given in fig-7(a) and -7(b). Obtained result also supports the the criteria to divide the dataset into 3 component.



(a) Major 2 components obtained by PCA and K-Means on that major 2 components

(b) Major 3 components obtained by PCA and K-Means on that major 2 components

Figure 7: PCA and K-Means of all the fuels

#### 4.2. Analysis of classified sub-dataset

Main dataset is divide into the three sub-dataset using silhouette criteria and K-Means clustering algorithm. By analysing the sub-dataset, obtained observations are discussed here.

##### 4.2.1. Subdataset-I

It contains 331 data points. Mostly short and middle level alkanes at high pressure and temperature are part of this dataset. Certain data points at low temperate and high pressure are also part this set i.e. propane at 950K and 21.4 atm.

Fuel		Temperature (K)	Pressure (atm)	Fuel Mole Fraction (%)	Oxygen Mole Fraction (%)	Equivalence Ratio	Data Points
Propane	max	1841	67.8	4	20	5	169
	min	950	1.12	0.15	0.75	0.5	
Butane	max	1761	5.5	2	13	2	55
	min	1230	1.03	0.5	3.25	0.5	
Pentane	max	1533	3.75	0.5	4	1	15
	min	1261	1.62	0.25	4	0.5	
Hexane	max	1475	3.6	0.42	4	1	16
	min	1237	1.67	0.21	4	0.5	
Heptane	max	1676	16.72	1.874	20.6	2	66
	min	1048	1.14	0.36	2.2	0.5	
Octane	max	1455	3.81	0.32	4	1	10
	min	1265	1.87	0.32	4	1	

Table 3: Summary of sub-data:II with label 0 using K-Means algorithm

#### 4.2.2. Subdataset-II

It contains 92 data points. it almost covers all type of alkanes(by length). Mainly data points are separated by low fuel and oxygen mole fraction percentage.

Fuel		Temperature (K)	Pressure (atm)	Fuel Mole Fraction (%)	Oxygen Mole Fraction (%)	Equivalence Ratio	Data Points
Propane	max	1687	2.24	0.05	0.25	1	5
	min	1505	2.13	0.05	0.25	1	
Butane	max	1761	2.16	0.05	0.325	1	3
	min	1531	2.02	0.05	0.325	1	
Heptane	max	1784	15.81	0.2	2.2	1	24
	min	1229	1.61	0.03	0.33	0.5	
Octane	max	1434	2.19	0.16	4	0.5	5
	min	1252	1.91	0.16	4	0.5	
Decane	max	1706	2.28	0.2	3.1	1.2	20
	min	1327	1.22	0.03	0.3875	0.64	
Dodecane	max	1657	15.72	0.05146	1.911	1	31
	min	1252	2.07	0.0371	0.731	0.5	
Hexadecane	max	1333	6.77	0.0497	1	1.22	4
	min	1170	4.44	0.0371	1	0.76	

Table 4: Summary of sub-data:I with label 1 using K-Means algorithm

#### 4.2.3. Subdataset-III

It contains 92 data points. Form intermediate to long chain alkanes at slightly lower temperature and high pressure along with high fuel and oxygen mole fractions are part of these dataset. Mainly it contains long chain alkanes with high fuel and oxygen content compared to other dataset.

Fuel		Temperature (K)	Pressure (atm)	Fuel Mole Fraction (%)	Oxygen Mole Fraction (%)	Equivalence Ratio	Data Points
Heptane	max	1115	60.6	1.874	20.6	1	16
	min	806	18.13	1.874	20.6	1	
Nonane	max	1301	41.76	0.4	4	2	27
	min	1051	13.52	0.2	4	0.5	
Decane	max	1173	5.15	2.567	21	1.89	5
	min	1081	4.56	1.44	21	1.06	
Dodecane	max	1422	33.7	2.138	21.0	1.88	131
	min	727	4	0.0558	2.786	0.05	
Hexadecane	max	1355	6.46	0.1832	4	1.22	14
	min	1159	1.71	0.0909	4	0.56	

Table 5: Summary of sub-data:III with label 2 using K-Means algorithm

### 5. Synthetic data generation and sampling technique:

To attain the objective, sufficient data points are necessary. From table- 3, 4, 5, it is clear that dataset is imbalance in terms of number of points for certain fuel. Coefficients obtained using multi-linear regression of such dataset will directly bias towards the more available information.

To reduce the bias in the result, different sampling strategy are utilized. This problem can be solved by two ways: [40]

1. Under-Sampling : Down sizing the largest dataset to minority size.
2. Over-Sampling : Expanding minority dataset size to majority dataset size

In under-sampling, negligence of data affect the correlation as less points may only cover limited range of physical condition. Due to drop of information, correlation may become detail specific. In over-sampling, expansion of data points is done by replicating the same data points which may ignore important information while sampling and may cause over-fitting. For IDT correlation rather than replicating samples uncertainty information is used to for over-sampling.

From table-1 it is clear that every fuel data-point contains certain range of error in measurement of temperature and pressure. Experimental error are mainly of two kinds:

1. Random Error : Occurs due to unknown and predictable changes
2. Systematic Error : Occurs due to measuring instrument and system handling

Error in measurement of shock temperature and pressure is uncontrolled phenomena which causes random error. Random error generally follows the Gaussian distribution. Using this information, more data points are generated using Multivariate Gaussian Normal Distribution.

$$p(x : \mu, \sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} \cdot |\Sigma|^{\frac{1}{2}}} \cdot \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \quad (23)$$

where,

$\Sigma$  is covariance matrix between parameter.

$\mu$  is vector of mean values or reported value of parameter

$n$  is number of parameters

#### 5.1. Algorithm to generate synthetic data points

---

##### Algorithm 1 Algorithm to generate random samples using uncertainty

---

```

1: NUF = Number of Unique Fuel in the dataset
2: MNDP = Maximum Number of Data Points in largest fuel set
3: DPSF = total number of Data Points in Selected Fuel
4: TNDP = Total Number of Data Points in Extended dataset
5: DPG = number of Data Points to be Generated from each sample
6:
7: TNDP = NUF * MNDP
8: for i = 1, 2, 3, ..., NUF in fuel data points do                                ▷ Select Unique Fuel
9:     DPG = TNDP / DPSF                                                         ▷ P2points from each sample
10:    for j = 1, 2, ..., DPSF do
11:        DATA_POINT_GENERATOR(data_point, DPG)                                ▷ Append to dataframe
12:    return Extended Dataframe                                                  ▷ Contains all fuel data
13: procedure DATA_POINT_GENERATOR(data_point, DPG)
14:      $\mu = [T\_means, P\_mean, \tau_{mean}]$                                        ▷ Reported values in the dataset
15:      $\Sigma = \begin{bmatrix} \sigma_T^2 & 0 & 0 \\ 0 & \sigma_P^2 & 0 \\ 0 & 0 & \sigma_\tau^2 \end{bmatrix}$                                ▷ From error values
16:     Generate 2000 points Using Multivariate Normal Gaussian  $\mathcal{N} \sim \mathcal{N}(\mu, \Sigma)$ 
17:     Random sampling by DPG count                                              ▷ to equi-size the dataset
18:     return Extended Dataframe                                                ▷ extended from single data point

```

---

To generate more data points using error values, reported value of temperature and pressure is used as mean values and given error range considered as standard error. Standard error values

taken as standard deviation to maximize the range of uncertainty. For single point n is considered as 1.

$$\sigma_{SD} = \frac{\sigma_{error}}{\sqrt{n}} \approx \sigma_{error} \quad (24)$$

Out of 9 dimensions, only 3 dimensions are considered to generate data points as other parameters are constant. Constant parameters were copied as it is in the data-frame as only temperature, pressure and IDT carries uncertainty. Reported value pressure and temperature is used for standard deviation. For IDT uncertainty varies maximum  $\pm 30\%$  but to be within bound,  $\pm 20\%$  uncertainty is used. Using multivariate Gaussian distribution along with relevant  $\Sigma$  and  $\mu$ , 2000 random points are generated to cover maximum possible occurrence. From 2000 random data points, sampling is done in such way that it will normalize the size of data points for all fuel. Procedure is mentioned in algorithm-1. To generate 2000 data points `np.random.multivariate_normal(mean, cov, 2000).T` function used and for sampling `random.choices(data, k = data_generation_count)` is utilized.

Once extended equi-sized data-frame is generated it data should be divided into 3 sub-dataset using K-means algorithm and then individual data-frame is transferred to Multi-linear regression module.

## 6. Multiple Regression : OLS estimator and hypothesis testing

Multi linear regression model is given by.

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \epsilon_i \quad (25)$$

where,

$Y_i = i^{th}$  observation of dependent variable

$X_{1i}, X_{2i}, X_{3i}, \dots, X_{ki}$  are independent observation of the k regression

$\epsilon_i$  is the error term which is not covered by independent variables

Using ordinary least square, estimation of coefficient is done using observed values.

$$\hat{Y}_i = b_0 + b_1 X_{1i} + b_2 X_{2i} + \dots + b_k X_{ki} \quad (26)$$

where,

$b_0, b_1, b_2, \dots, b_k$  are the estimated value of the  $\beta_0, \beta_1, \beta_2, \dots, \beta_k$  and coefficient  $b_0, b_1, b_2, \dots, b_k$  are obtained by solving linear system of equation.

$$\begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & X_{13} & \dots & X_{1k} \\ 1 & X_{21} & X_{22} & X_{23} & \dots & X_{2k} \\ 1 & X_{31} & X_{32} & X_{33} & \dots & X_{3k} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & X_{n3} & \dots & X_{nk} \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_k \end{bmatrix}$$

For IDT dataset, dependent and independent variables are mentioned in formulation-13. In IDT equation-16, it is assumed that ignition delay depends on all the parameters and obtained from over-determinant system ( $n > k$ ) of equation using least-square solution. Objective function



is defined as,

$$\hat{b} = \underset{b}{\operatorname{argmin}} \quad ||y - Xb||_2 \quad (27)$$

where,  $\hat{b}$  is coefficient for best fit rather than the actual solution.

$$\begin{aligned} \hat{b} &= \underset{b}{\operatorname{argmin}} \quad ||y - Xb||_2 \\ &= \underset{b}{\operatorname{argmin}} \sum_{i=1}^n |y_i - \sum_{j=1}^k X_{ij}b_j| \\ &= \underset{b}{\operatorname{argmin}} \quad (y - Xb)^T (y - Xb) \\ &= \underset{b}{\operatorname{argmin}} \quad (y^T - b^T X^T)(y - Xb) \\ &= \underset{b}{\operatorname{argmin}} \quad y^T y - y^T Xb - b^T X^T y + b^T X^T Xb \end{aligned} \quad (28)$$

Differentiating with respect to  $b^T$ ,

$$\begin{aligned} 0 &= -2X^T y + 2X^T Xb \\ 2X^T y &= 2X^T Xb \end{aligned} \quad (29)$$

$$b = (X^T X)^{-1} X^T y \quad (30)$$

By solving equation-30 all coefficient are obtained. While modelling of IDT, it was assumed that IDT depends on all the parameters. But significance of each independent variable on dependent variable varies and has to be verified by hypothesis testing.

### 6.1. Hypothesis testing

Null Hypothesis ( $H_0$ ): There is no significant relationship between independent variable and dependent variable or in another words, by adding an independent variable to model will not casue any significant improvement on the dependent variable. Means  $b_1 = b_2 = b_3 = \dots = b_n = 0$ .

Hypothesis testing done using t-test by taking confidence interval of 95%. To obtain coefficient and other other statical parameter, statmodel [41] library package is used. In which, t-values is obtained by taking ratio of coefficient and standard error. P-values (probability values) are obtained from t-values using two-tail test.

If obtained p-values is less than  $\alpha = 0.05$  (5%), which means value lies in the confidence interval and obtained result or independent parameters is statically significant to dependent parameter. In 95% confidence interval or  $p < 0.05$  generally null hypothesis is rejected. To select statistically significant independent variable backward elimination method is used [42]. Independent variable  $X_i$  were removed till removal does not cause drastic decrease in  $R^2$  value.

After eliminating all non-significant independent variable, obtained final correlation and result are discussed further.

---

**Algorithm 2** Backward elimination

---

- 1: Use all independent variable to obtain correlation
  - 2: find coefficient and p values of independent variable
  - 3: **while** for any  $X_i : p > 0.05$  **do**
  - 4:     Remove  $X_i$  from dataset
  - 5:     Obtain coefficient and p values
- 

**7. Result and Discussion:**

Most important part this study is K-means clustering. Location of cluster centroid is important for unknown fuel to select the cluster category. Movement of cluster centroids with variation in data point are mentioned in table-6. From the table it is observed that cluster centroids are very sensitive to the number of data points and fuels. But as number of data points, variety fuel, physical conditions increases cluster centroid converges to specific location and movement in centroid become less. For individual fuel, cluster centroid varies significantly but as combination of fuel increases it converges to location as in case of all fuel.

Fuels vs Axis		Tempeprature	Pressure	$C_p - C_H$	$C_p - C_S$	$C_S - C_S$	$C_S - C_H$	Fuel(%)	Oxygen(%)
All	Centroid-0	1.57837308	1.51675986	1.24959195	0.41653065	0.26505381	0.94663827	-0.23781285	1.84398964
	Centroid-1	1.68012488	0.93518365	1.12307542	0.37435847	1.22535434	2.82506715	-2.88900159	-0.19240077
	Centroid-2	1.38326654	2.36889142	1.51647253	0.50549084	2.11119433	4.7278795	-0.6030923	2.56675347
Propane, Dodecane, Hexane, Butane, Heptane, Decane	Centroid-0	1.57605773	1.57148828	1.25323871	0.41774624	0.23445554	0.88665731	-0.16905458	1.88849474
	Centroid-1	1.6944866	0.91018368	1.1061271	0.36870903	1.14345847	2.65562598	-2.92404759	-0.29666696
	Centroid-2	1.35803663	2.42946412	1.55515259	0.5183842	2.16787013	4.85412445	-0.37914573	2.88516679
Propane, Dodecane, Heptane	Centroid-0	1.5444353	1.96307245	1.29529591	0.4317653	0.25076071	0.93328672	-0.11458897	1.97025854
	Centroid-1	1.35250581	2.35543233	1.56435254	0.52145085	2.26644565	5.05434215	-0.41749343	2.92323109
	Centroid-2	1.67874594	1.13744756	1.12470066	0.37490022	1.1954392	2.76577862	-2.95741867	-0.28904548
Propane	Centroid-0	1.59652546	1.25275661	1.22966736	0.40988912	0.	0.40988912	0.02766885	1.87766756
	Centroid-1	1.42465876	3.62448678	1.4507666	0.48358887	0.	0.48358887	-0.127854	1.92640905
	Centroid-2	1.75337899	0.77994782	1.0400941	0.34669803	0.	0.34669803	1.75337899	0.77994782
Heptane	Centroid-0	1.6708165	1.22938158	1.1350263	0.3783421	0.7566842	1.89171049	-3.18913949	-0.65261478
	Centroid-1	1.65276957	0.5790902	1.15044562	0.38348187	0.76696375	1.91740937	-0.57119456	1.83443768
	Centroid-2	1.33539693	3.39234322	1.59378677	0.53126226	1.06252451	2.65631128	0.62807518	3.02529108
Dodecane	Centroid-0	1.49689856	2.6228982	1.34469596	0.44823199	2.01704393	4.48231985	-2.529175	1.33872879
	Centroid-1	1.3509045	2.26022207	1.56610377	0.52203459	2.34915565	5.22034588	-0.29324417	3.04132966
	Centroid-2	1.70214704	0.80548421	1.0960673	0.36535577	1.64410095	3.65355766	-3.07720187	-0.2855489

Table 6: Location of 3- centroids obtained using K-Means for various of fuel/fuel-combinations

From each sub-dataset 80% data is used as training set and remaining 20% is used as testing set for performing multiple regression and hypothesis testing. Obtained coefficient after complete procedure-22 is mentioned in table-7. For multiple regression greedy binary tree algorithm is implemented, in which if minimum  $R^2=0.85$  training accuracy is not obtained then algorithm will start dividing the data into two parts based on mean temperature value of dataset. All sub-dataset satisfies that cut-off criteria and obtained coefficient are mentioned in table-7.

	Intercept	Tempeprature	Pressure	$C_p - C_H$	$C_p - C_S$	$C_S - C_S$	$C_S - C_H$	Fuel(%)	Oxygen(%)	Training $R^2$	Test $R^2$
Cluster -0	37.1733	-16.1974	-0.5486	-0.9622	-0.3207	0.119	0	0.9398	-1.6978	0.95186	0.95405
	33.6319	-14.8109	-0.5521	0	0	0.1346	0	0.9329	-1.6929	0.95336	0.95254
Cluster -1	38.5577	-17.2162	-0.4416	-1.7754	-0.5918	0.2792	0	0.67	-1.1244	0.96630	0.97684
	32.4963	-14.9362	-0.4456	0	0	0.2809	0	0.6651	-1.1229	0.97353	0.96695
Cluster -2	101.8355	-42.4864	-1.0402	-19.2446	-6.4149	2.5527	-1.3095	-0.1454	-0.3402	0.88752	0.91019
	102.0124	-42.5188	-1.0379	-19.3288	-6.4429	2.5669	-1.3091	-0.1297	-0.348	0.91015	0.88578

Table 7: coefficient and accuracy obtained for 3 sub-dataset using training set of all available fuel data points

Due to generation of data points from uncertainty and sampling from generated points, causes variation in result. Such variation is acceptable as obtained accuracy for both test and train set

falls within acceptable range. To observe the variation more precisely and find out the effect of sampling on coefficient, results were obtained by running the simulation 4400 times. Variation in the result is plotted by histograms:8-18. In the results, distribution of coefficients follows the nearly Gaussian distribution. for coefficient distribution deviates due to classification error. But as overall accuracy falls within acceptable range such results are considerable.

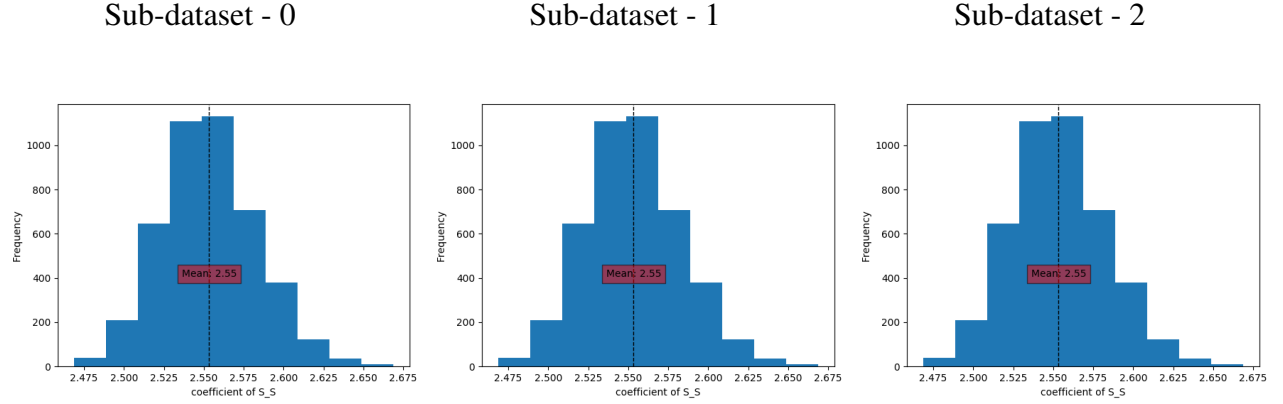


Figure 8: Secondary-Secondary Carbon bond coefficient of sub-dataset

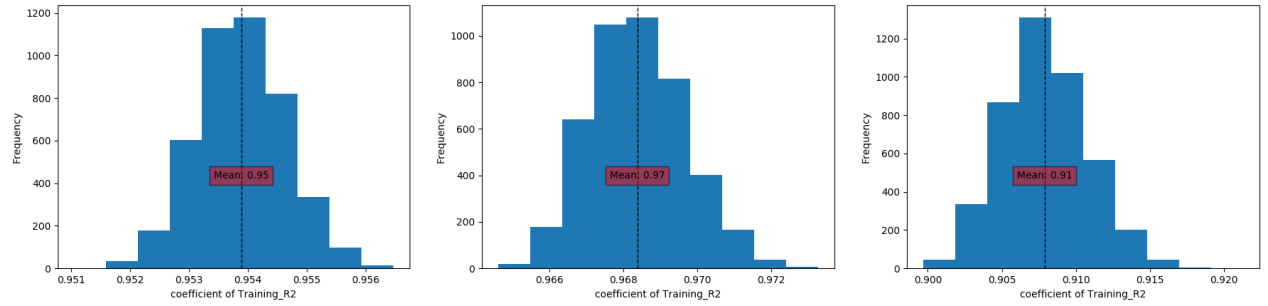


Figure 9: Training accuracy of sub-dataset

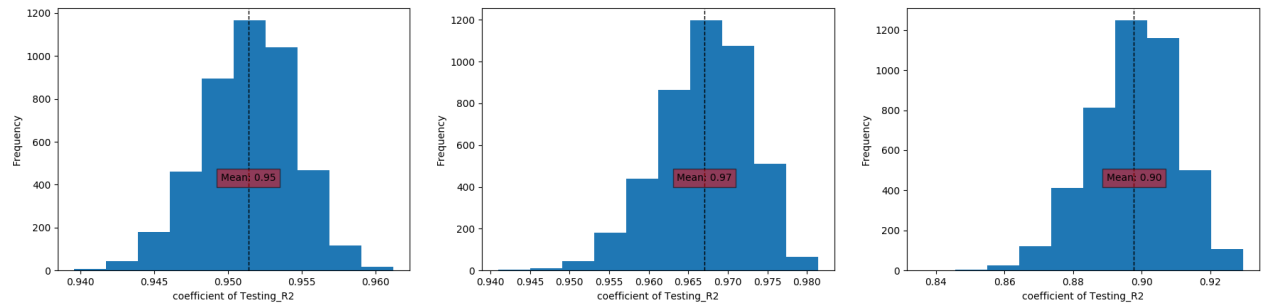
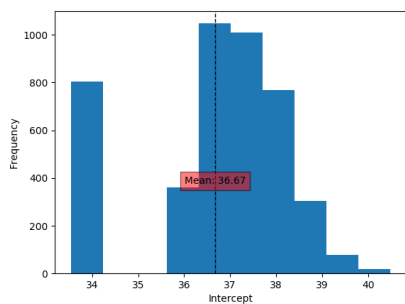
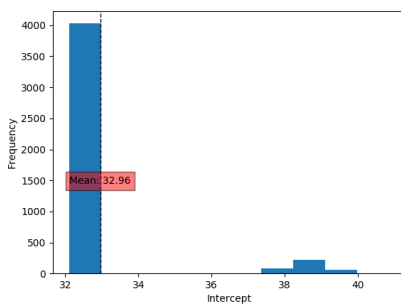


Figure 10: Testing accuracy of sub-dataset

Sub-dataset - 0



Sub-dataset - 1



Sub-dataset - 2

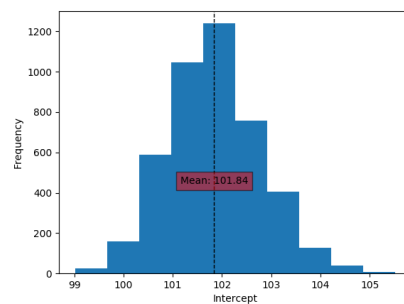


Figure 11: Intercept of sub-dataset

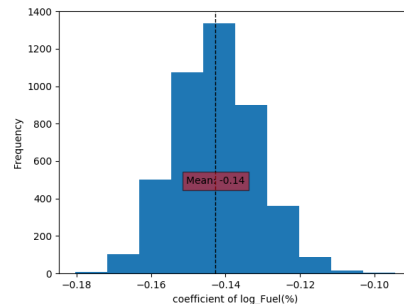
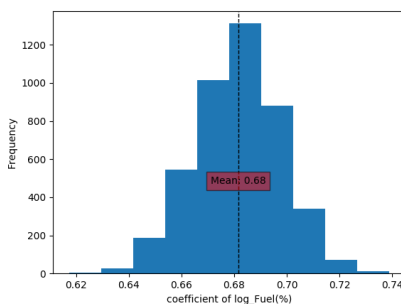
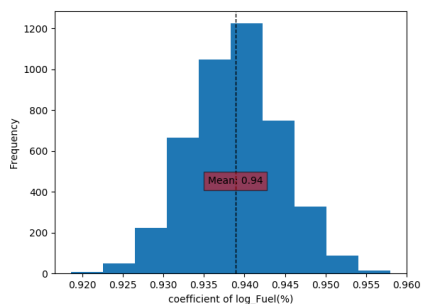


Figure 12: Fuel mole fraction term coefficient of sub-dataset

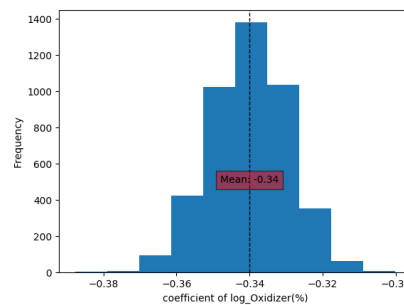
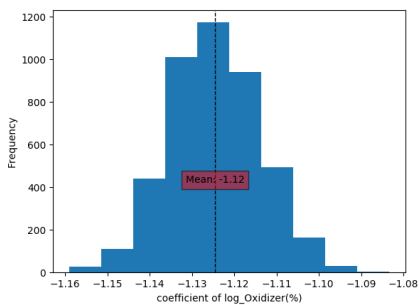
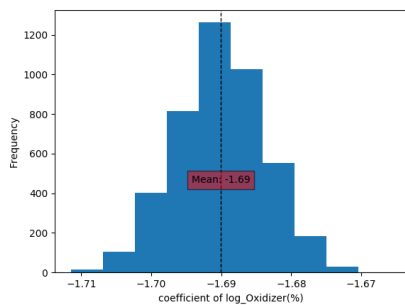


Figure 13: Oxygen term coefficient of sub-dataset

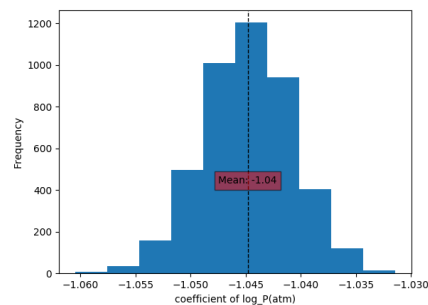
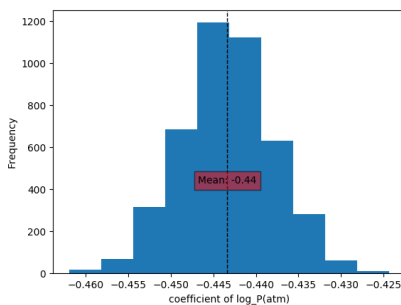
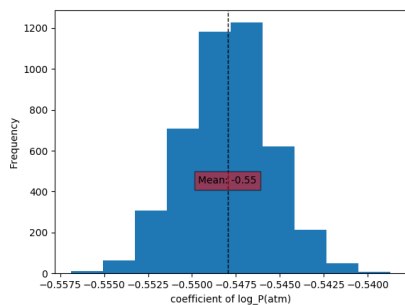
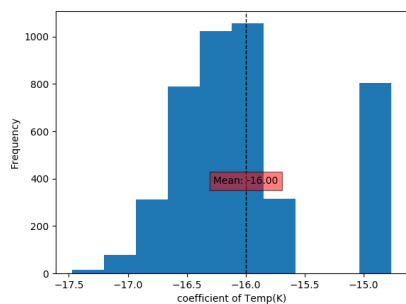
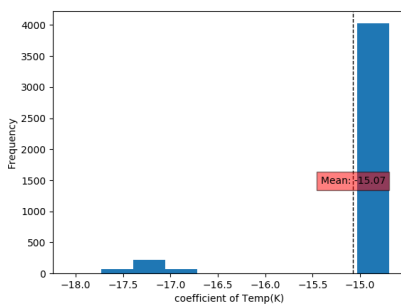


Figure 14: Pressure term coefficient of sub-dataset

Sub-dataset - 0



Sub-dataset - 1



Sub-dataset - 2

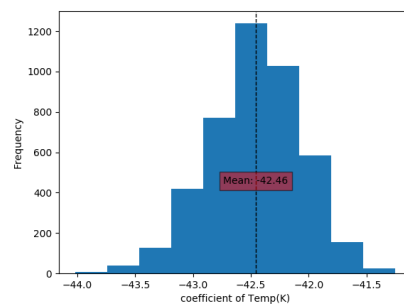


Figure 15: Temperature term coefficient of sub-dataset

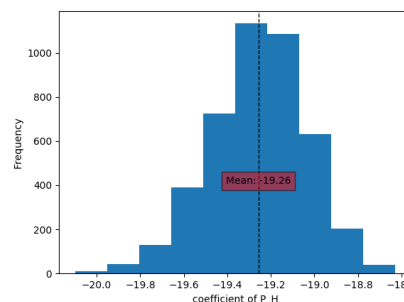
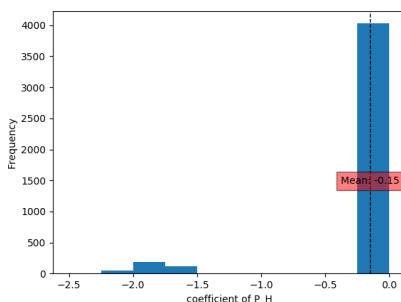
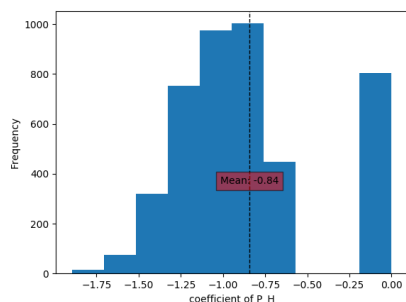


Figure 16: Primary Carbon and Hydrogen bond coefficient of sub-dataset

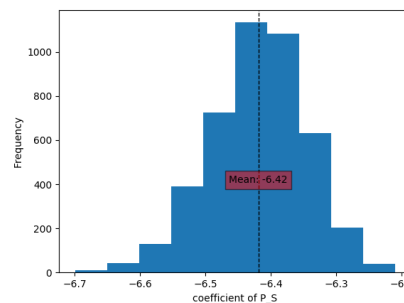
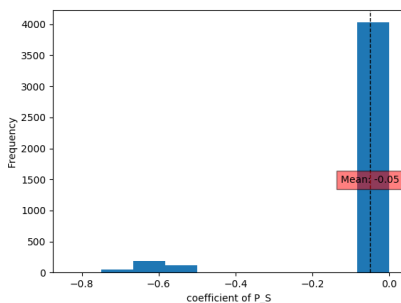
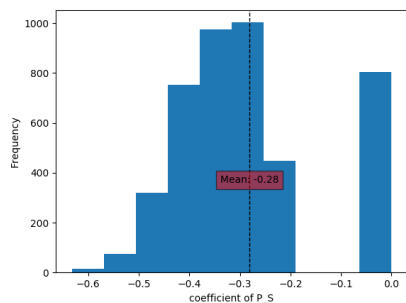


Figure 17: Primary-Secondary Carbon bond coefficient of sub-dataset

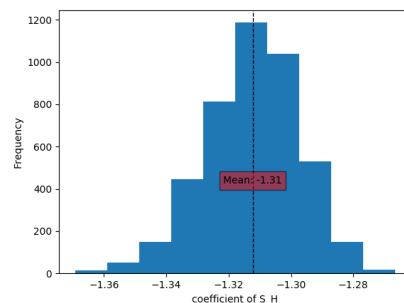
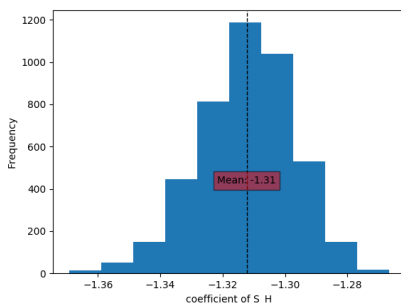
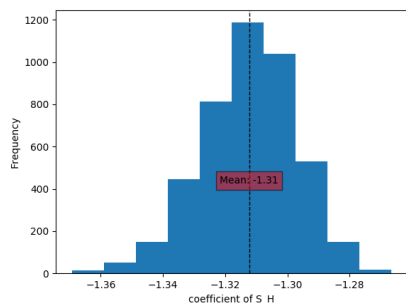


Figure 18: Secondary Carbon and Hydrogen bond coefficient of sub-dataset

Formulation of Ignition delay time is acquired using mean value of such distribution. After 4200 repetitive result mean value of coefficient converges to specific value as given in table-8. IDT correlation obtained from converged coefficient value- 8 is given as below. Associated value of centroids are mentioned in table-6 for 'all fuel' case:

- For centroid-0:

$$\tau = e^{36.68} \cdot \left(\frac{T}{T_0}\right)^{-15.99} \cdot \left(\frac{P}{P_0}\right)^{-0.55} \cdot X_{Fuel}^{0.94} \cdot X_{O_2}^{-1.69} \cdot e^{(C_P H)^{-0.84}} \cdot e^{(C_P C_S)^{-0.28}} \cdot e^{(C_S C_S)^{-0.13}} \quad (31)$$

- For centroid-1:

$$\tau = e^{32.96} \cdot \left(\frac{T}{T_0}\right)^{-15.08} \cdot \left(\frac{P}{P_0}\right)^{-0.44} \cdot X_{Fuel}^{0.68} \cdot X_{O_2}^{-1.12} \cdot e^{(C_P H)^{-0.16}} \cdot e^{(C_P C_S)^{-0.05}} \cdot e^{(C_S C_S)^{-0.29}} \quad (32)$$

- For centroid-2:

$$\tau = e^{101.85} \cdot \left(\frac{T}{T_0}\right)^{-42.46} \cdot \left(\frac{P}{P_0}\right)^{-1.04} \cdot X_{Fuel}^{-0.14} \cdot X_{O_2}^{-0.34} \cdot e^{(C_P H)^{-19.26}} \cdot e^{(C_P C_S)^{-6.42}} \cdot e^{(C_S C_S)^{-2.55}} \cdot e^{(C_S H)^{-1.31}} \quad (33)$$

	Number of Simulation	Constant	P_H	P_S	S_H	S_S	Temp(K)	log_Fuel(%)	log_Oxidizer(%)	log_P(atm)	Testing_R2	Training_R2
Cluster-0	500	36.67	-0.84	-0.28	0	0.13	-16	0.94	-1.69	-0.55	0.95	0.95
	1000	36.72	-0.86	-0.29	0	0.13	-16.02	0.94	-1.69	-0.55	0.95	0.95
	1500	36.67	-0.84	-0.28	0	0.13	-15.99	0.94	-1.69	-0.55	0.95	0.95
	2000	36.66	-0.84	-0.28	0	0.13	-15.99	0.94	-1.69	-0.55	0.95	0.95
	2500	36.65	-0.84	-0.28	0	0.13	-15.99	0.94	-1.69	-0.55	0.95	0.95
	3000	36.65	-0.84	-0.28	0	0.13	-15.99	0.94	-1.69	-0.55	0.95	0.95
	3500	36.67	-0.84	-0.28	0	0.13	-16	0.94	-1.69	-0.55	0.95	0.95
	4000	36.68	-0.85	-0.28	0	0.13	-15.99	0.94	-1.69	-0.55	0.95	0.95
	4200	36.68	-0.84	-0.28	0	0.13	-15.99	0.94	-1.69	-0.55	0.95	0.95
	4395	36.68	-0.84	-0.28	0	0.13	-15.99	0.94	-1.69	-0.55	0.95	0.95
Cluster-1	500	33.16	-0.21	-0.07	0	0.28	-15.15	0.68	-1.12	-0.44	0.97	0.97
	1000	33.04	-0.18	-0.05	0	0.29	-15.1	0.68	-1.12	-0.44	0.97	0.97
	1500	33.01	-0.17	-0.06	0	0.29	-15.09	0.68	-1.12	-0.44	0.97	0.97
	2000	32.99	-0.16	-0.05	0	0.29	-15.08	0.68	-1.12	-0.44	0.97	0.97
	2500	32.98	-0.16	-0.05	0	0.29	-15.08	0.68	-1.12	-0.44	0.97	0.97
	3000	32.97	-0.16	-0.05	0	0.29	-15.08	0.68	-1.12	-0.44	0.97	0.97
	3500	32.97	-0.16	-0.05	0	0.29	-15.08	0.68	-1.12	-0.44	0.97	0.97
	4000	32.96	-0.16	-0.05	0	0.29	-15.08	0.68	-1.12	-0.44	0.97	0.97
	4200	32.96	-0.16	-0.05	0	0.29	-15.08	0.68	-1.12	-0.44	0.97	0.97
	4395	32.96	-0.16	-0.05	0	0.29	-15.08	0.68	-1.12	-0.44	0.97	0.97
Cluster-2	500	101.83	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	1000	101.85	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	1500	101.86	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	2000	101.84	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	2500	101.84	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	3000	101.86	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	3500	101.85	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	4000	101.84	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	4200	101.84	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91
	4395	101.85	-19.26	-6.42	-1.31	2.55	-42.46	-0.14	-0.34	-1.04	0.9	0.91

Table 8: Coefficient values obtained after given number of simulations for different sub-dataset which shows convergence of coefficient values

For verification of correlation and frame work, smaller set dataset of heptane and hexadecane from, main dataset used to verify the result. All fuel data points used for training the model heptane and hexadecane are also part of it. After complete procedure mentioned as in 22, observed result shows excellent match with predicted value. For cluster-0,1 relative error is bounded between

0.8%-3% and 0.4%-5% respectively [19](#). For cluster-2 error varied from 0.09% to 16% which is high as obtained coefficient of determination is also low.

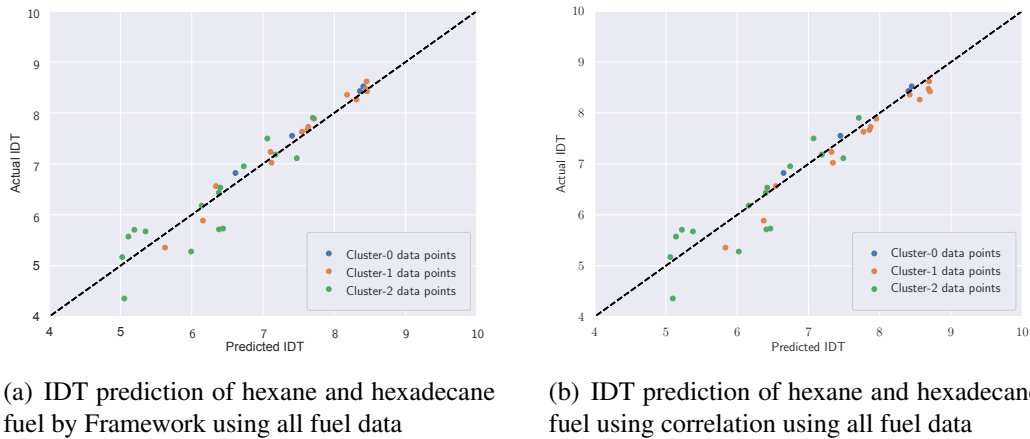


Figure 19: Prediction of IDT of hexane and hexadecane using correlation and Framework

Result of IDT from correlation shows slightly more deviation. But it follows the same trend as of framework result. Good prediction of hexane and hexadecane is expected as correlation already contains those fuel detail. The goal is present study is to predict IDT for new unknown fuel.

To check prediction of unknown straight chain alkane fuel, Hexadecane and Hexane is used as unknown fuel which would not be part of training and testing set. Model is trained using remaining all fuel. To attain main objective and check behaviour of framework as well as correlational result are mentioned in figure-[20](#).

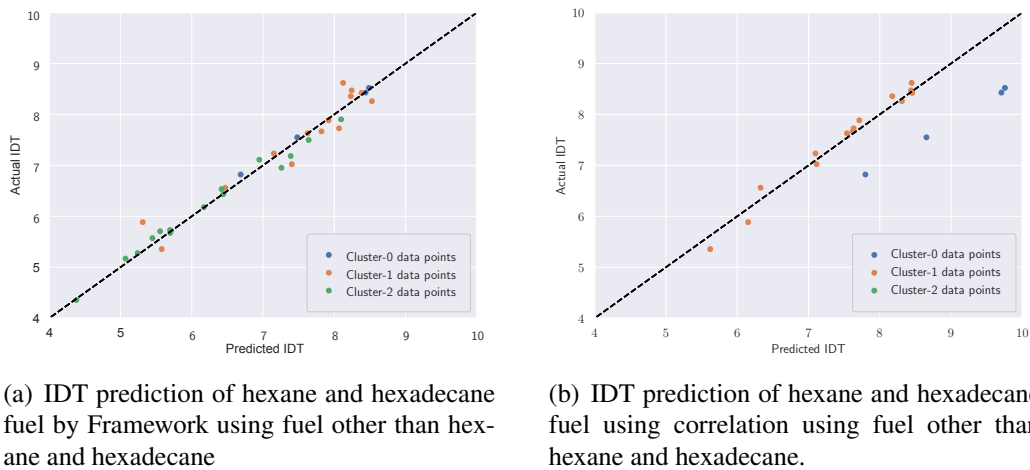
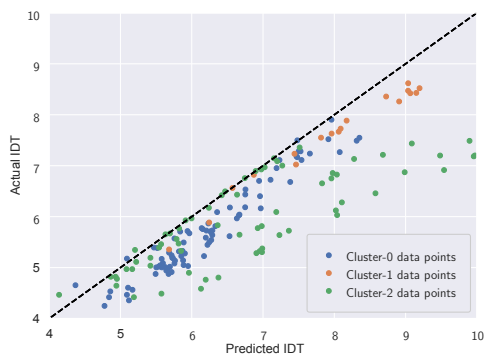
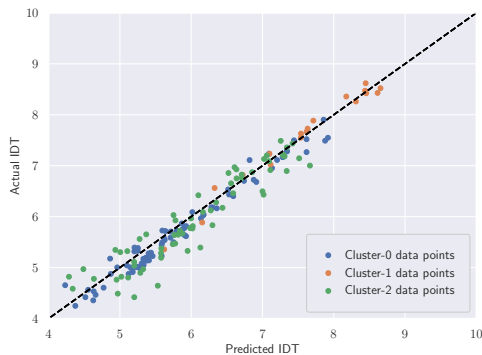


Figure 20: Prediction of IDT of hexane and hexadecane using correlation and Framework

Result of figure-[20](#) shows that prediction using framework again gives excellent fit but result obtained using correlation incurred considerable error in prediction. Error in correlation result might be caused due to the movement of centroids which shows that Accuracy of correlational depends also on centroid value.



(a) IDT prediction of all fuel other than propane, heptane and dodecane by Framework using propane, heptane and dodecane as learning set



(b) IDT prediction of all fuel other than propane, heptane and dodecane using correlation using propane, heptane and dodecane as learning set.

Figure 21: Prediction of IDT of hexane and hexadecane using correlation and Framework

To verify extreme case, propane, heptane and dodecane were used to train the model as it almost covers wide range of condition and different length of alkanes. Rest of all fuel were used as unknown fuel to test the prediction. Observed result was quite unexpected, which is shown in fig-21. Obtained result from framework is quite scattered specifically for cluster-2. Whereas result obtained using correlation [31,32,33](#) shows excellent agreement of predicted result with experimental result. Centroids were obtained from propane, heptane and dodecane. Such pair of centroid and correlation is important for prediction IDT of new fuel.



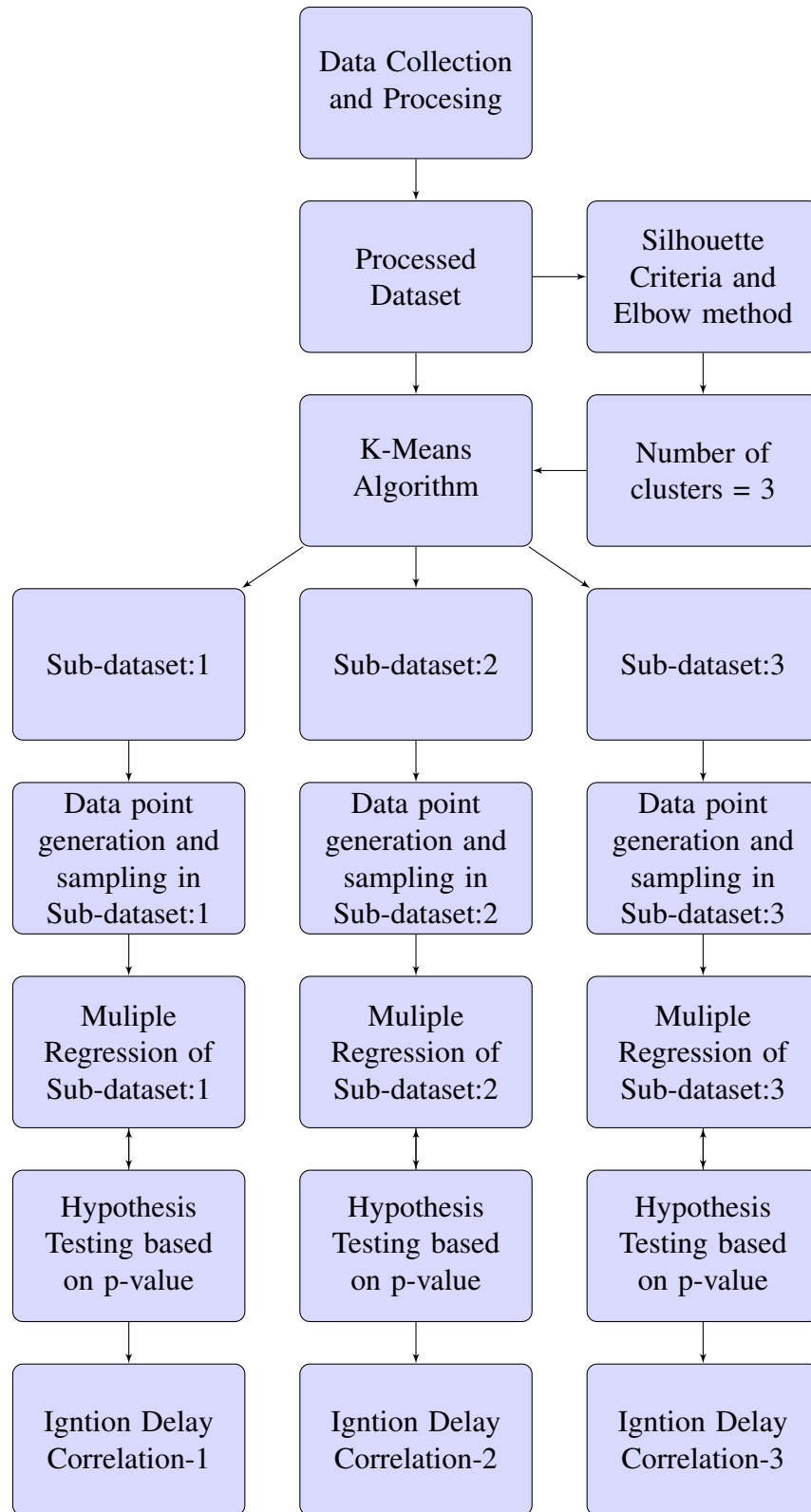


Figure 22: Flowchart of Ignition Delay Time Prediction Framework

## References:

- [1] T. Lu and C. K. Law, "Toward accommodating realistic fuel chemistry in large-scale computations," *Progress in Energy and Combustion Science*, vol. 35, no. 2, pp. 192–215, 2009.
- [2] S. S. Goldsborough, "A chemical kinetically based ignition delay correlation for iso-octane covering a wide range of conditions including the ntc region,"
- [3] D. C. Horning, "A study of the high-temperature auto-ignition and thermal decomposition of hydrocarbons," vol. Report No. TSD-135, 2001.
- [4] F. Khaled and A. Farooq, "On the universality of ignition delay times of distillate fuels at high temperatures: A statistical approach," *Combustion and Flame*, vol. 210, pp. 145–158, 2019.
- [5] E. V. Anslyn and D. A. Dougherty, *Modern physical organic chemistry*. University science books, 2006.
- [6] S. J. Blanksby and G. B. Ellison, "Bond dissociation energies of organic molecules," *Accounts of chemical research*, vol. 36, no. 4, pp. 255–263, 2003.
- [7] M. Frenklach, H. Wang, M. Goldenberg, G. Smith, D. Golden, C. Bowman, R. Hanson, W. Gardiner, and V. Lissianski, "Gri-mech—an optimized detailed chemical reaction mechanism for methane combustion," *Gas Research Institute Report No. GRI-95/0058*, 1995.
- [8] M. Röhrig, E. L. Petersen, D. F. Davidson, R. K. Hanson, and C. T. Bowman, "Measurement of the rate coefficient of the reaction  $\text{CH} + \text{O}_2 \rightarrow \text{products}$  in the temperature range 2200 to 2600 K," *International journal of chemical kinetics*, vol. 29, no. 10, pp. 781–789, 1997.
- [9] J. Zhang, E. Hu, Z. Zhang, L. Pan, and Z. Huang, "Comparative study on ignition delay times of C1–C4 alkanes," *Energy & Fuels*, vol. 27, no. 6, pp. 3480–3487, 2013.
- [10] E. Hu, Y. Chen, Z. Zhang, X. Li, Y. Cheng, and Z. Huang, "Experimental study on ethane ignition delay times and evaluation of chemical kinetic models," *Energy & Fuels*, vol. 29, no. 7, pp. 4557–4566, 2015.
- [11] D. C. Horning, D. Davidson, and R. Hanson, *A study of the high-temperature autoignition and thermal decomposition of hydrocarbons*. PhD thesis, Stanford University Stanford, California, 2001.
- [12] D. C. Horning, D. Davidson, and R. Hanson, "Study of the high-temperature autoignition of n-alkane/O<sub>2</sub>/Ar mixtures," *Journal of Propulsion and Power*, vol. 18, no. 2, pp. 363–371, 2002.
- [13] D. Davidson, J. Herbon, D. Horning, and R. Hanson, "OH concentration time histories in n-alkane oxidation," *International journal of chemical kinetics*, vol. 33, no. 12, pp. 775–783, 2001.
- [14] K. Y. Lam, *Shock tube measurements of oxygenated fuel combustion using laser absorption spectroscopy*. PhD thesis, Stanford University, 2013.

- [15] K.-Y. Lam, Z. Hong, D. Davidson, and R. Hanson, "Shock tube ignition delay time measurements in propane/o<sub>2</sub>/argon mixtures at near-constant-volume conditions," *Proceedings of the Combustion Institute*, vol. 33, no. 1, pp. 251–258, 2011.
- [16] S. M. Burke, U. Burke, R. Mc Donagh, O. Mathieu, I. Osorio, C. Keesee, A. Morones, E. L. Petersen, W. Wang, T. A. DeVerter, *et al.*, "An experimental and modeling study of propene oxidation. part 2: Ignition delay time and flame speed measurements," *Combustion and Flame*, vol. 162, no. 2, pp. 296–314, 2015.
- [17] D. Davidson, S. Ranganath, K.-Y. Lam, M. Liaw, and Z. Hong, "Ignition delay time measurements of normal alkanes and simple oxygenates," *Journal of Propulsion and Power*, vol. 26, no. 2, pp. 280–287, 2010.
- [18] B. Gauthier, D. F. Davidson, and R. K. Hanson, "Shock tube determination of ignition delay times in full-blend and surrogate fuel mixtures," *Combustion and Flame*, vol. 139, no. 4, pp. 300–311, 2004.
- [19] D. Davidson, M. Oehlschlaeger, and R. Hanson, "Methyl concentration time-histories during iso-octane and n-heptane oxidation and pyrolysis," *Proceedings of the Combustion Institute*, vol. 31, no. 1, pp. 321–328, 2007.
- [20] S. S. Vasu, D. F. Davidson, and R. K. Hanson, "Oh time-histories during oxidation of n-heptane and methylcyclohexane at high pressures and temperatures," *Combustion and Flame*, vol. 156, no. 4, pp. 736–749, 2009.
- [21] D. Davidson, Z. Hong, G. Pilla, A. Farooq, R. Cook, and R. Hanson, "Multi-species time-history measurements during n-heptane oxidation behind reflected shock waves," *Combustion and flame*, vol. 157, no. 10, pp. 1899–1905, 2010.
- [22] D. R. Haylett, *The development and application of aerosol shock tube methods for the study of low-vapor-pressure fuels*. Stanford University, 2011.
- [23] D. Davidson, D. Haylett, and R. Hanson, "Development of an aerosol shock tube for kinetic studies of low-vapor-pressure fuels," *Combustion and Flame*, vol. 155, no. 1-2, pp. 108–117, 2008.
- [24] D. R. Haylett, D. F. Davidson, and R. K. Hanson, "Ignition delay times of low-vapor-pressure fuels measured using an aerosol shock tube," *Combustion and Flame*, vol. 159, no. 2, pp. 552–561, 2012.
- [25] D. Haylett, D. Davidson, and R. Hanson, "Second-generation aerosol shock tube: an improved design," *Shock Waves*, vol. 22, no. 6, pp. 483–493, 2012.
- [26] D. Haylett, D. Davidson, and R. Hanson, "A second-generation aerosol shock tube for combustion research," in *48th AIAA Aerospace Sciences Meeting Including the New Horizons Forum and Aerospace Exposition*, p. 196, 2010.

- [27] D. Jackson, D. Davidson, and R. Hanson, "Application of an aerosol shock tube for the kinetic studies of n-dodecane/nano-aluminum slurries," in *44th AIAA/ASME/SAE/ASEE joint propulsion conference & exhibit*, p. 4767, 2008.
- [28] S. S. Vasu, D. F. Davidson, Z. Hong, V. Vasudevan, and R. K. Hanson, "n-dodecane oxidation at high-pressures: Measurements of ignition delay times and oh concentration time-histories," *Proceedings of the Combustion Institute*, vol. 32, no. 1, pp. 173–180, 2009.
- [29] S. S. Vasu, *Measurements of ignition times, OH time-histories, and reaction rates in jet fuel and surrogate oxidation systems*. PhD thesis, PhD Thesis, Stanford University, California, United States, 2010.
- [30] D. Haylett, D. Davidson, R. Cook, Z. Hong, W. Ren, S. Pyun, and R. Hanson, "Multi-species time-history measurements during n-hexadecane oxidation behind reflected shock waves," *Proceedings of the Combustion Institute*, vol. 34, no. 1, pp. 369–376, 2013.
- [31] D. Haylett, R. Cook, D. Davidson, and R. Hanson, "Oh and c<sub>2</sub>h<sub>4</sub> species time-histories during hexadecane and diesel ignition behind reflected shock waves," *Proceedings of the Combustion Institute*, vol. 33, no. 1, pp. 167–173, 2011.
- [32] D. Weininger, "Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules," *Journal of chemical information and computer sciences*, vol. 28, no. 1, pp. 31–36, 1988.
- [33] G. Landrum, "Rdkit: Open-source cheminformatics,"
- [34] W. Ji, P. Zhao, T. He, X. He, A. Farooq, and C. K. Law, "On the controlling mechanism of the upper turnover states in the ntc regime," *Combustion and Flame*, vol. 164, pp. 294–302, 2016.
- [35] Z. Zhao, Z. Chen, and S. Chen, "Correlations for the ignition delay times of hydrogen/air mixtures," *Chinese science bulletin*, vol. 56, no. 2, pp. 215–221, 2011.
- [36] G. Ogbuabor and F. Ugwoke, "Clustering algorithm for a healthcare dataset using silhouette score value," *International Journal of Computer Science & Information Technology*, vol. 10, no. 2, pp. 27–37, 2018.
- [37] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [38] P. J. Rousseeuw, "Silhouettes: a graphical aid to the interpretation and validation of cluster analysis," *Journal of computational and applied mathematics*, vol. 20, pp. 53–65, 1987.
- [39] R. C. de Amorim and C. Hennig, "Recovering the number of clusters in data sets with noise features using feature rescaling factors," *Information Sciences*, vol. 324, pp. 126–145, 2015.

- 406 [40] N. V. Chawla, “Data mining for imbalanced datasets: An overview,” in *Data mining and*  
407 *knowledge discovery handbook*, pp. 875–886, Springer, 2009.
- 408 [41] S. Seabold and J. Perktold, “statsmodels: Econometric and statistical modeling with python,”  
409 in *9th Python in Science Conference*, 2010.
- 410 [42] J. H. McDonald, *Handbook of biological statistics*, vol. 2. 2009.