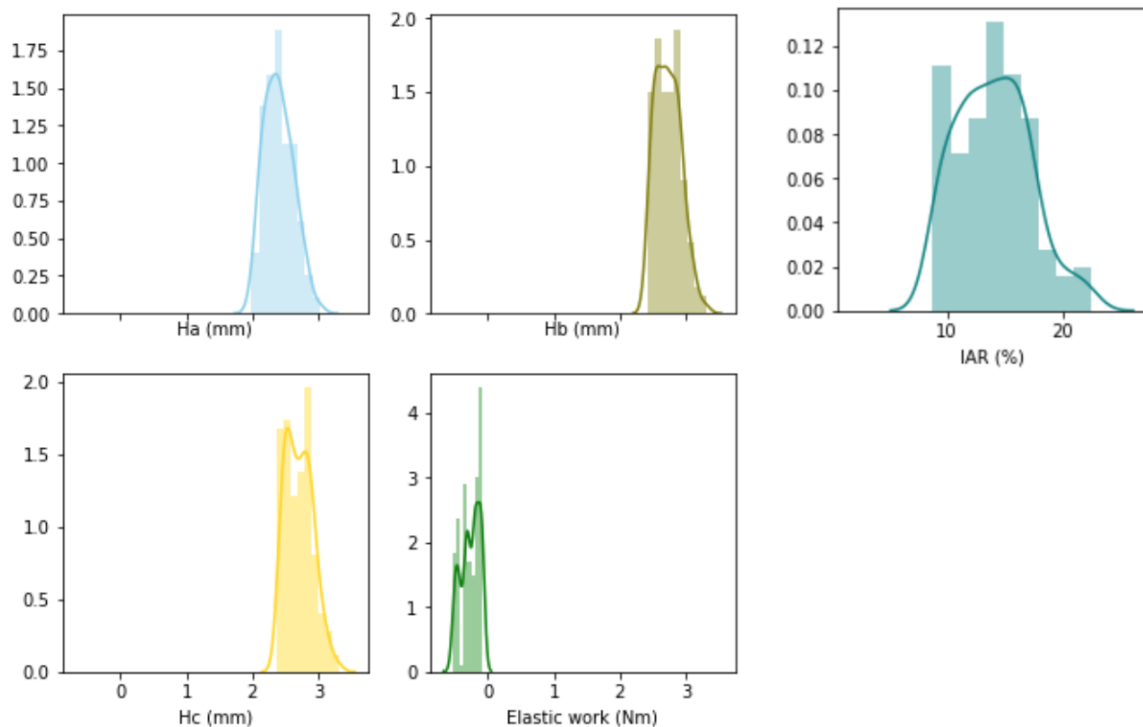# 1. Introduction

Many issues in pharmaceutical manufacturing are direct consequence of tablet expansion which is caused by elastic behaviour of materials. To ensure quality of a product and facilitate industrial scale-up and manufacturing processes Quality by design principles must be employed. Prediction of elastic properties of granules if very important in pre-formulation and formulation studies. Therefore, prediction of elastic properties is very important in wet granulation scale-up process. In this experiment I set up few hypothesis testing to understand the data in better way and then, I try to evaluate different predictive models for the same data.

# 2. Statistical Data Analysis

## 2.1 Distribution



From the distribution of Response variables (Ha, Hb, Hc, IAR, ElasticWork), it was suspected that the sample is not normally distributed. To confirm, **Shapiro-Wilk test** had been carried out and the result mimic the same way with our initial assumption (p-value < 0.05).

**Shapiro-Wilk test Result:**

```
> shapiro.test(Data$Ha..mm.)

        Shapiro-Wilk normality test

data:  Data$Ha..mm.
W = 0.97894, p-value = 0.01249

> shapiro.test(Data$Hb..mm.)

        Shapiro-Wilk normality test

data:  Data$Hb..mm.
W = 0.96802, p-value = 0.0007038

> shapiro.test(Data$Hc..mm.)

        Shapiro-Wilk normality test

data:  Data$Hc..mm.
W = 0.96294, p-value = 0.0002082

> shapiro.test(Data$Elastic.work..Nm.)

        Shapiro-Wilk normality test

data:  Data$Elastic.work..Nm.
W = 0.90891, p-value = 1.211e-08

> shapiro.test(Data$IAR....)

        Shapiro-Wilk normality test

data:  Data$IAR....
W = 0.96999, p-value = 0.001151
```

## 2.2 Data Bucketing

As the response variables are not normally distributed, I may use Non Parametric Approach (Kruskal-Wallis test) to find relation between run (laboratory, commercial, pilot) and the response variables to find whether they should be treated differently.

It was found that test ran in different environment(laboratory, commercial, pilot) has made an different impact on the elastic property of the tablet, so I divided Run into buckets such as Laboratory, Commercial and pilot for our further analysis.

```
> kruskal.test(Data$Ha..mm.~Data$Run)

        Kruskal-Wallis rank sum test

data:  Data$Ha..mm. by Data$Run
Kruskal-Wallis chi-squared = 23.322, df = 2, p-value = 8.623e-06

> kruskal.test(Data$Hb..mm.~Data$Run)

        Kruskal-Wallis rank sum test

data:  Data$Hb..mm. by Data$Run
Kruskal-Wallis chi-squared = 19.575, df = 2, p-value = 5.614e-05

> kruskal.test(Data$Hc..mm.~Data$Run)

        Kruskal-Wallis rank sum test

data:  Data$Hc..mm. by Data$Run
Kruskal-Wallis chi-squared = 37.083, df = 2, p-value = 8.864e-09

> kruskal.test(Data$IAR....~Data$Run)

        Kruskal-Wallis rank sum test

data:  Data$IAR.... by Data$Run
Kruskal-Wallis chi-squared = 43.769, df = 2, p-value = 3.132e-10

> kruskal.test(Data$Elastic.work..Nm.~Data$Run)

        Kruskal-Wallis rank sum test

data:  Data$Elastic.work..Nm. by Data$Run
Kruskal-Wallis chi-squared = 5.8947, df = 2, p-value = 0.05248
```
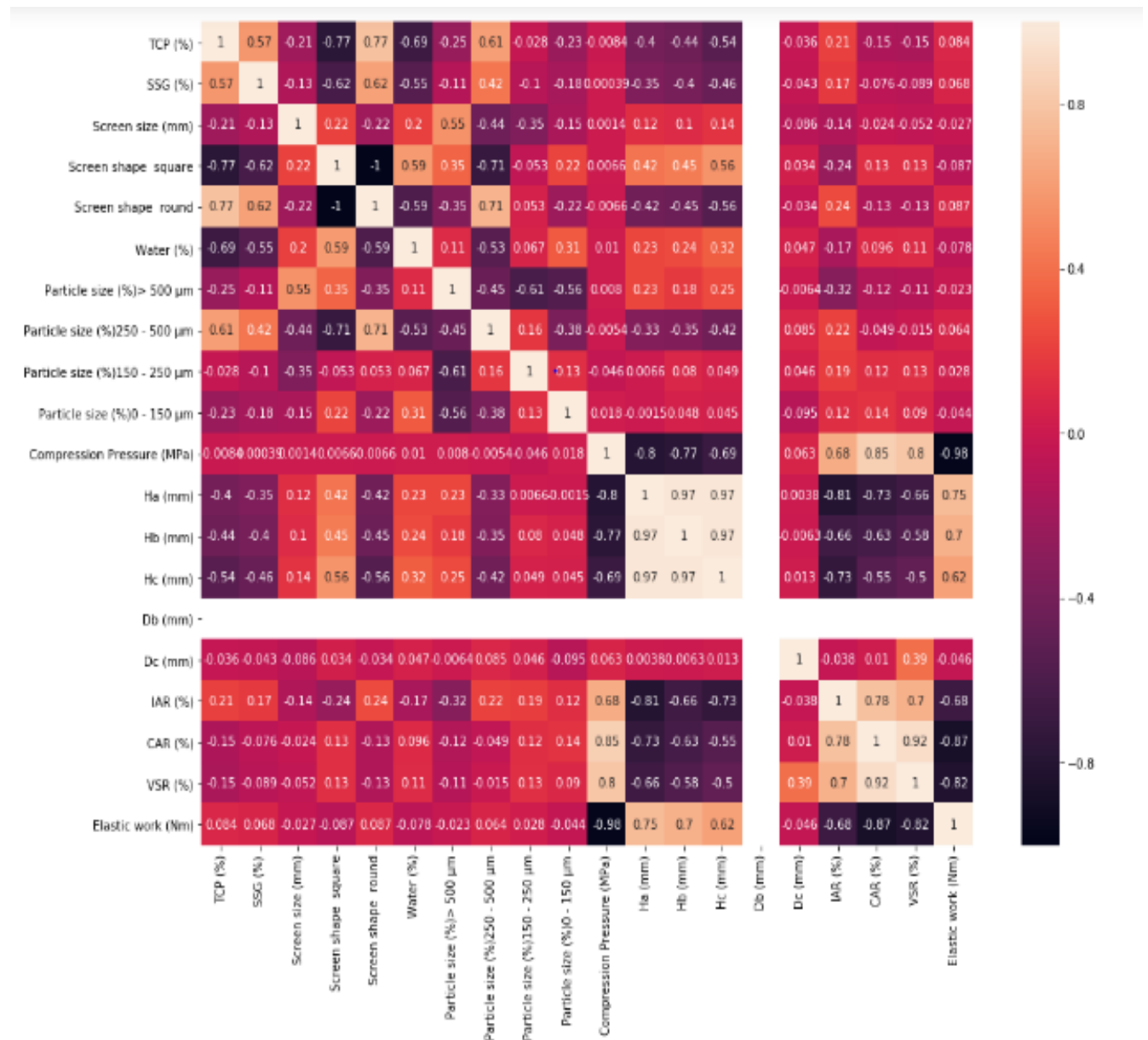
## 2.2 Relation between Dependent and Independent Variables

Before modelling, it was required to analyse the impact of independent variables on dependent variables (Ha, Hb, Hc, IAR, ElasticWork). As part of that process I am trying to get the correlation matrix between these variables.

## Correlation Matrix:



Correlation Matrix of the Data was used for feature selection

# 3. Models:

I have carried out a regression analysis on the input and output variables as given below:

**Input Variables:**

1. TCP (%)
2. SSG (%)
3. Screen size (mm)
4. Screen shape  square
5. Screen shape  round
6. Water (%)
7. Particle size (%)> 500 µm
8. Particle size (%)250 - 500 µm
9. Particle size (%)150 - 250 µm
10. Particle size (%)0 - 150 µm
11. Compression Pressure (MPa)

Output Variable:

1. Ha (mm)
2. Hb (mm)
3. Hc (mm)
4. IAR (%)
5. Elastic work (Nm)

I have evaluated the following regression based algorithms for predicting response variables.

- Step-wise Liner Regression
- Linear Regression with Ridge
- Linear Regression with Lasso
- Linear Regression with ElasticNet

I have a concern of overfitting for this small dataset; so I thought to use regularisation variants of regression algorithms for my analysis.

### 3.1 Model for Response Ha:

#### 3.1.1 Laboratory scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
Ha

Model with Best Score/ Low MSE:
**0.8425 / 0.03 (Linear Regression)**

#### 3.1.2 Pilot scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
Ha

Model with Best Score/ Low MSE:
**0.8881 / 0.01 (Linear Regression with Ridge)**

#### 3.1.3 Commercial scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
Ha

Model with Best Score/ Low MSE:

**0.8508 / 0.01 (Linear Regression)**

### 3.2 Model for Response Hb:

#### 3.2.1 Laboratory scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Screen shape  round','Particle size (%)150 - 250 µm'

Output Variable:
Hb

Model with Best Score/ Low MSE:
**0.7873 / 0.03 (Linear Regression)**

### 3.2.2 Pilot scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Screen shape  round'

Output Variable:
Hb

Model with Best Score/ Low MSE:
**0.8843 / 0.004 (Linear Regression)**

### 3.2.3 Commercial scale

Input Variables Considered (Best Score/Low MSE):
'Compression Pressure (MPa)'

Output Variable:
Hb

Model with Best Score/ Low MSE:
**0.8736 / 0.003 (Linear Regression with ElasticNet)**

## 3.3 Model for Response Hc:

### 3.3.1 Laboratory scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
Hc

Model with Best Score/ Low MSE:
**0.7786 / 0.03 (Linear Regression)**

### 3.3.2 Pilot scale

Input Variables Considered (Best Score/Low MSE):

'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
Hc

Model with Best Score/ Low MSE:
**0.8461 / 0.01 (Linear Regression with Ridge)**

### 3.3.3 Commercial scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
Hc

Model with Best Score/ Low MSE:
**0.8534 / 0.004 (Linear Regression)**

## 3.4 Model for Response IAR:

### 3.4.1 Laboratory scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
IAR

Model with Best Score/ Low MSE:
**0.9052 / 1.47 (Linear Regression)**

### 3.4.2 Pilot scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','Compression Pressure (MPa)','Particle size (%)150 - 250 µm','Screen shape round'

Output Variable:
IAR

Model with Best Score/ Low MSE:
**0.9168 / 0.63 (Linear Regression with Ridge)**

### 3.4.3 Commercial scale

Input Variables Considered (Best Score/Low MSE):
'TCP (%)','SSG (%)','Particle size (%)> 500 μm','Compression Pressure (MPa)','Particle size (%)150 - 250 μm','Screen shape  round'

Output Variable:
IAR

Model with Best Score/ Low MSE:
**0.4556 / 7.55 (Linear Regression)**

## 3.5 Model for Response Elastic Work:

### 3.5.1 Laboratory scale

Input Variables Considered (Best Score/Low MSE):
'Compression Pressure (MPa)'

Output Variable:
Elastic Work

Model with Best Score/ Low MSE:
**0.9576 / 0.001 (Linear Regression)**

### 3.5.2 Pilot scale

Input Variables Considered (Best Score/Low MSE):
'Compression Pressure (MPa)'

Output Variable:
Elastic Work

Model with Best Score/ Low MSE:
**0.9858 / 0.0002 (Linear Regression)**

### 3.5.3 Commercial scale

Input Variables Considered (Best Score/Low MSE):
'Compression Pressure (MPa)'

Output Variable:
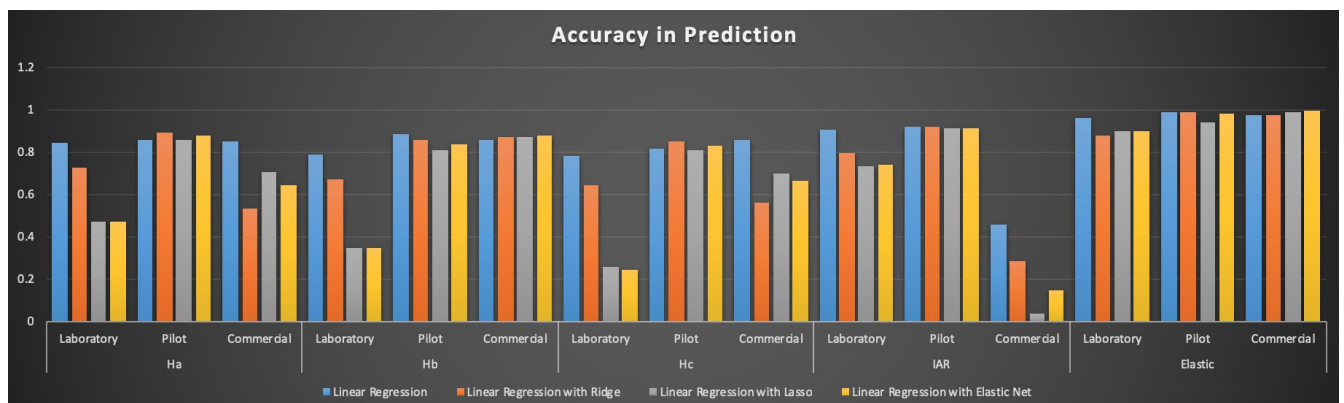Elastic Work

Model with Best Score/ Low MSE:
**0.9902 / 0.001 (Linear Regression with ElasticNet**

# 4. Results:

The table describes the Accuracy and Mean Squared Error (MSE) from all four different regression models built for four response variables. The best results are highlighted in the table.

| Response Variable | Scale | Linear Regression | | Linear Regression with Ridge | | Linear Regression with Lasso | | Linear Regression with Elastic | |
|---|---|---|---|---|---|---|---|---|---|
| | | Result | MSE | Result | MSE | Result | MSE | Result | MSE |
| Ha | Laboratory Scale | 0.8425 | 0.03 | 0.7257 | 0.01 | 0.4714 | 0.02 | 0.4684 | 0.02 |
| | Pilot Scale | 0.8562 | 0.01 | 0.8881 | 0.01 | 0.8571 | 0.01 | 0.8771 | 0.01 |
| | Commercial Scale | 0.8508 | 0.01 | 0.533 | 0.01 | 0.7013 | 0.01 | 0.6451 | 0.01 |
| Hb | Laboratory Scale | 0.7873 | 0.03 | 0.6683 | 0.01 | 0.3474 | 0.02 | 0.3484 | 0.02 |
| | Pilot Scale | 0.8843 | 0.004 | 0.8545 | 0.01 | 0.8082 | 0.01 | 0.8381 | 0.01 |
| | Commercial Scale | 0.8567 | 0.004 | 0.8699 | 0.003 | 0.8666 | 0.003 | 0.8736 | 0.003 |
| Hc | Laboratory Scale | 0.7786 | 0.03 | 0.6435 | 0.01 | 0.2534 | 0.03 | 0.2433 | 0.03 |
| | Pilot Scale | 0.8113 | 0.01 | 0.8461 | 0.01 | 0.8043 | 0.01 | 0.8309 | 0.01 |
| | Commercial Scale | 0.8534 | 0.004 | 0.5594 | 0.01 | 0.6984 | 0.004 | 0.6601 | 0.004 |
| IAR | Laboratory Scale | 0.9052 | 1.47 | 0.796 | 1.09 | 0.7316 | 1.43 | 0.7398 | 1.39 |
| | Pilot Scale | 0.915 | 0.86 | 0.9168 | 0.63 | 0.911 | 0.68 | 0.9109 | 0.68 |
| | Commercial Scale | 0.4556 | 7.55 | 0.2815 | 4.21 | 0.037 | 5.83 | 0.1468 | 4.99 |
| Elastic | Laboratory Scale | 0.9576 | 0.001 | 0.8763 | 0.002 | 0.9002 | 0.001 | 0.8999 | 0.001 |
| | Pilot Scale | 0.9858 | 0.0002 | 0.9882 | 0.0002 | 0.9398 | 0.001 | 0.9776 | 0.0004 |
| | Commercial Scale | 0.9717 | 0.0009 | 0.9721 | 0.0004 | 0.9863 | 0.0002 | 0.9902 | 0.001 |

## 4.1 Model Accuracy



It quite evident from the above graph that, for the response variable 'Elastic Work', almost all the models predicted with similar high accuracy. However, for other response variables, stepwise regression provided a better result.

## 4.2 Model MSE

From the Mean Squared Analysis, it is shown that, Commercial Model Error is generally less compared to other buckets like Laboratory or Pilot. Also, with regularisation effect the models are not highly overfitted. In few occasions, Ridge Regression models gave more reliable results than others.