# Analysis of Different Datasets on Model Performance: Without Feature Selection vs. With Feature Selection

Tarun Chintada
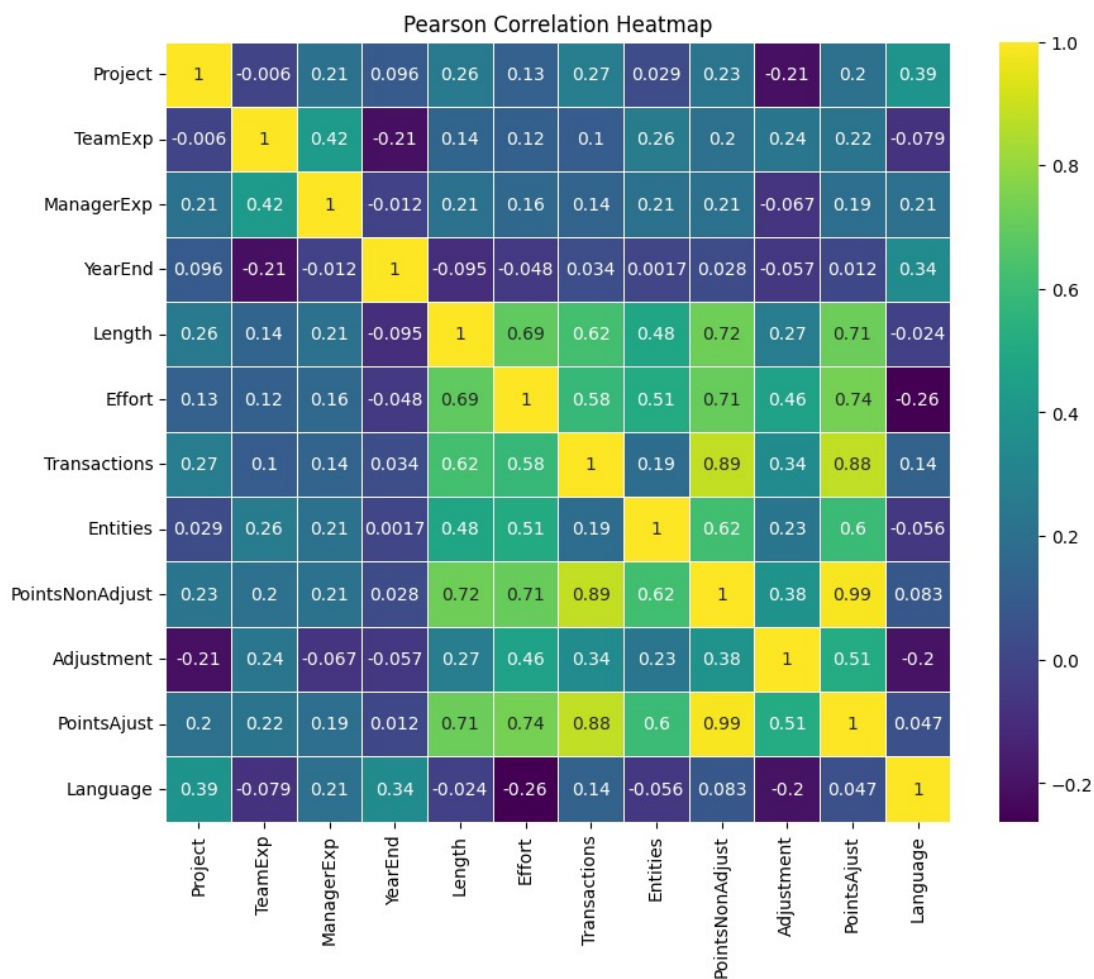
June 29, 2024

## Neural Networks(30-30) Regression

## 1. Desharnais Dataset

The Desharnais dataset analyzed consists of information from 81 software projects from a Canadian company. Each project has 12 attributes

**feature selection**



Pearson Correlation Heatmap

The selected features are Length ,Transactions,Entities,PointsNonAdjust ,PointsAjust which are highly correlated features whose correlation coefficient is greater than 0.5



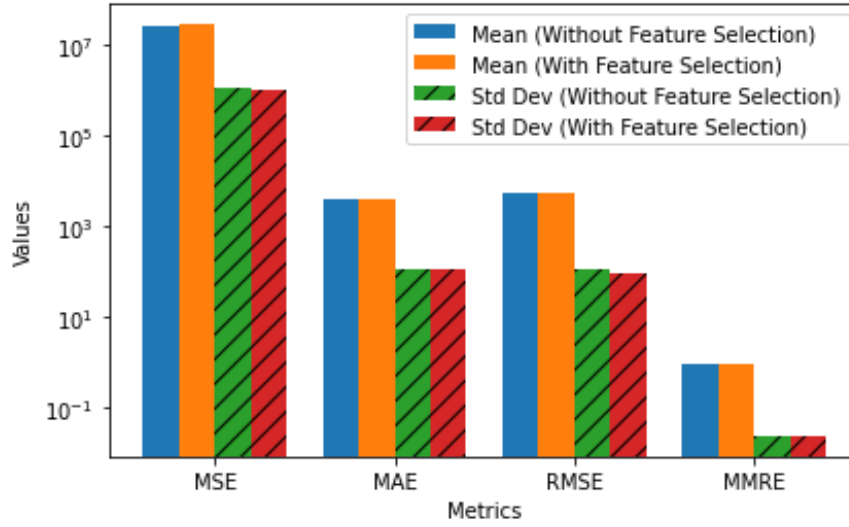Comparison of Mean and Std Dev For Desharnais Dataset With and out Feature Selction

**TABLE 1** .Comparison of the results Means and Standard Deviation (SD) of Deshrnais dataset

| Measures | Without Feature Selection | | With Feature Selection | |
|----------|------------|------------|------------|------------|
| | **Mean** | **SD** | **Mean** | **SD** |
| **MSE** | 27843741.3613 | 1199167.6066 | 28757662.3042 | 1021246.4407 |
| **MAE** | 4031.8336 | 111.1507 | 4013.2488 | 109.9385 |
| **RMSE** | 5275.4759 | 114.4352 | 5361.7675 | 95.4550 |
| **MMRE** | 0.8889 | 0.0245 | 0.8848 | 0.0242 |

## 2. albrecht Dataset

albrecht consist of 21 samples each sample has 9 attributes.

### feature selection



Pearson Correlation Heatmap

The selected features are Input,Output,Inquiry,File,RawFPcounts,AdjFP which are highly correlated features whose correlation coefficient is greater than 0.5



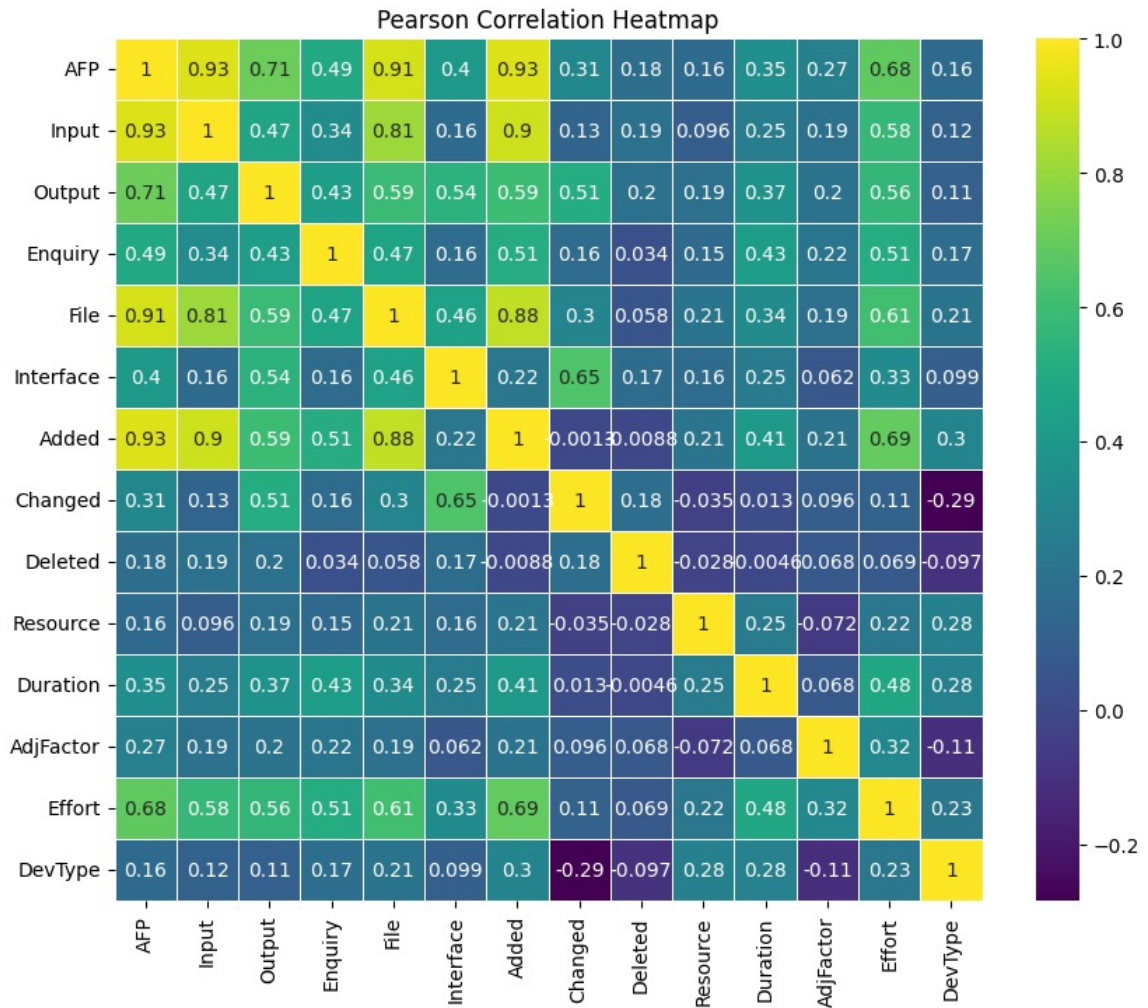Comparison of Mean and Std Dev For albreht Dataset With and out Feature Selction

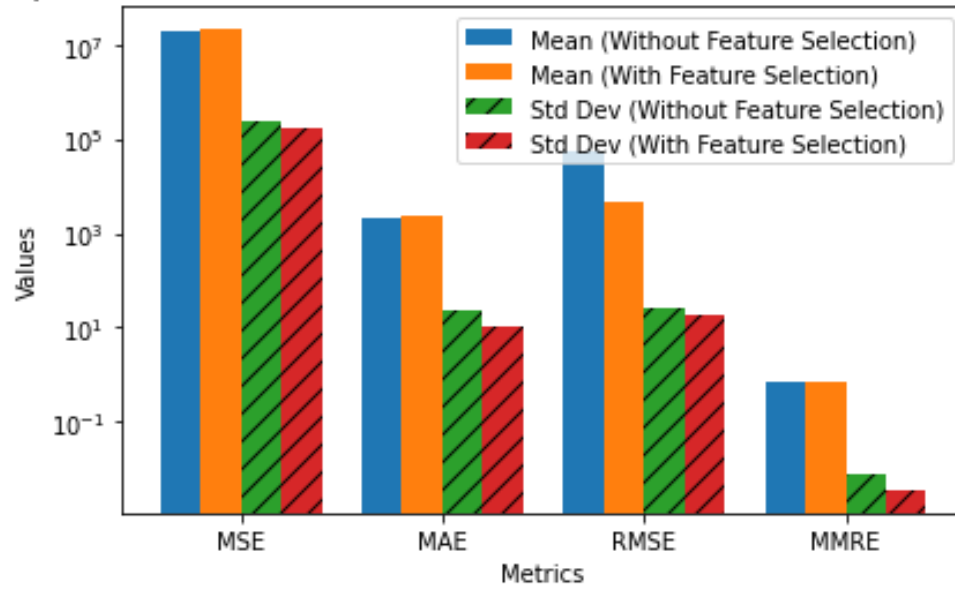| Measures | Without Feature Selection | | With Feature Selection | |
|----------|------|------|------|------|
| | Mean | SD | Mean | SD |
| MSE | 415.8218 | 227.5610 | 520.1206 | 291.2830 |
| MAE | 14.5934 | 3.6008 | 15.2890 | 4.0386 |
| RMSE | 19.6571 | 5.4241 | 21.9569 | 6.1656 |
| MMRE | 0.4179 | 0.1031 | 0.4378 | 0.1157 |

## 3. china Dataset

china consist of 499 samples each sample has 15 attributes.

**feature selection**



The selected features are AFP,Input,Output,Enquiry,File,Added which are highly correlated features whose correlation coefficient is greater than 0.5

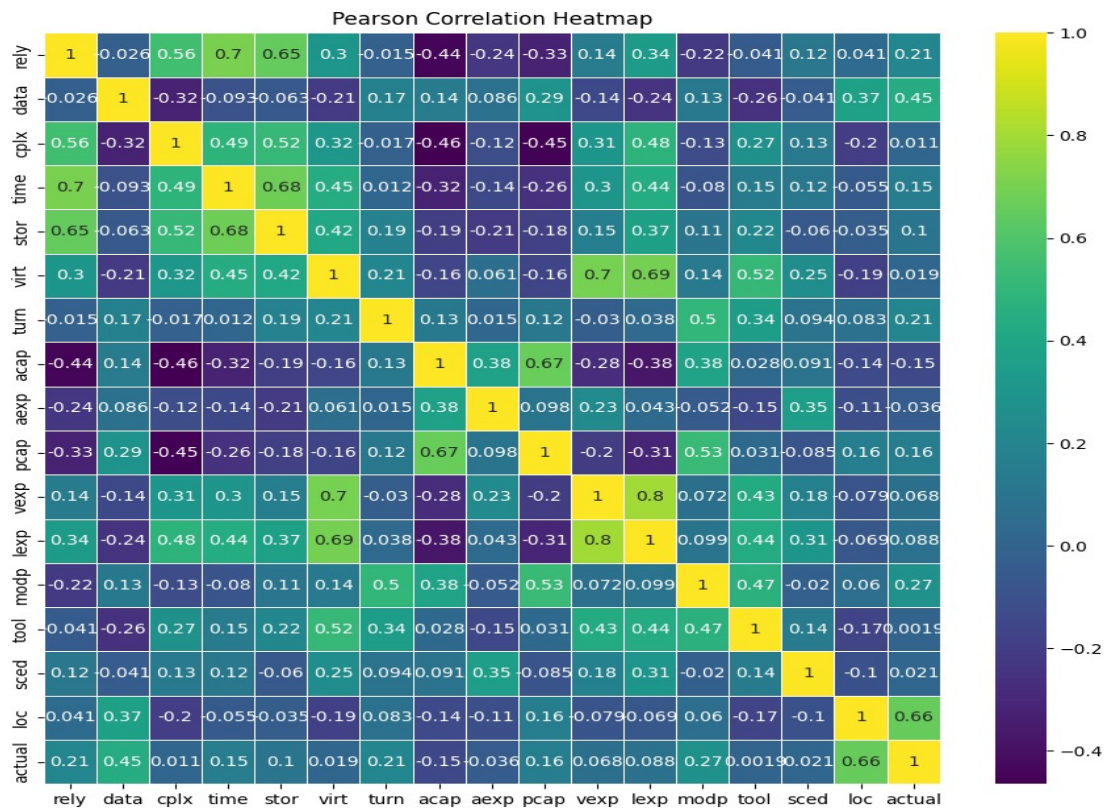Comparison of Mean and Std Dev For China Dataset With and out Feature Selction

| Measures | Without Feature Selection | | With Feature Selection | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| **MSE** | 18476460.7450 | 234433.2416 | 20838461.3713 | 171352.7200 |
| **MAE** | 2101.3670 | 24.1657 | 2314.0458 | 11.0766 |
| **RMSE** | 54298.3388 | 27.2824 | 4564.8778 | 18.7588 |
| **MMRE** | 0.6773 | 0.0078 | 0.7458 | 0.0242 |

# 4. cocomo81 Dataset

cocomo81 consist of 63 samples each sample has 17 attributes.

**feature selection**



Pearson Correlation Heatmap

The selected features are loc which are highly correlated features whose correlation coefficient is greater than 0.5



Comparison of Mean and Std Dev For cocomo81 Dataset With and out Feature Selction
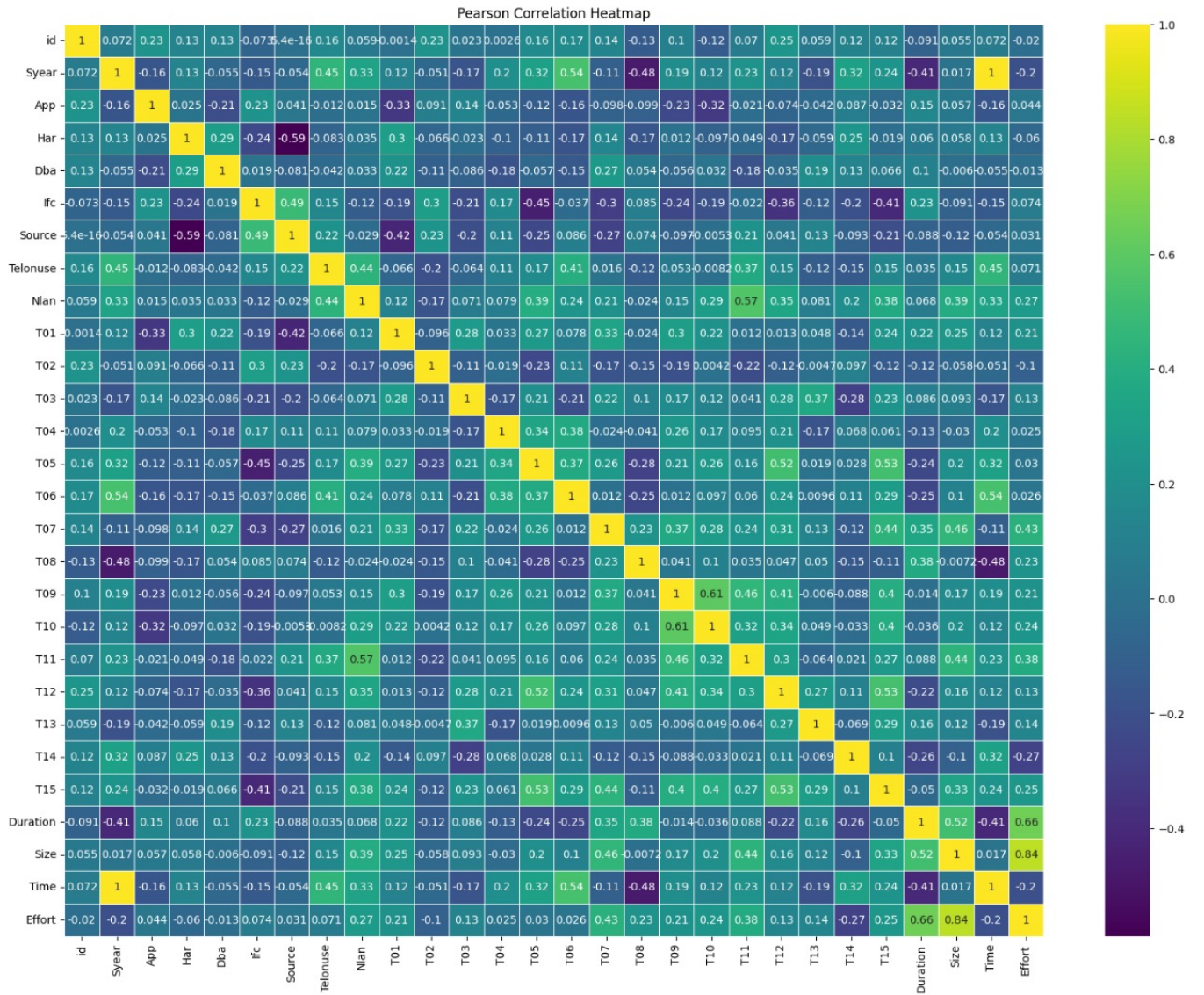
| Measures | Without Feature Selection | | With Feature Selection | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| MSE | 53921.8849 | 14041.7981 | 309291.2282 | 5368.2188 |
| MAE | 256.7712 | 15.4838 | 236.7179 | 4.7079 |
| RMSE | 503.7120 | 14.0038 | 556.1186 | 4.8267 |
| MMRE | 0.9826 | 0.0593 | 0.9059 | 0.0242 |

## 5.Maxwell Dataset

Maxwell consist of 62 samples each sample has 28 attributes.

## feature selection


Pearson Correlation Heatmap

The selected features are rely,data,turn,modp,loc which are highly correlated features whose correlation coefficient is greater than 0.2

Comparison of Mean and Std Dev For maxwell Dataset With and out Feature Selction

| Measures | Without Feature Selection | | With Feature Selection | |
|----------|---------------------------|--------------|------------------------|--------------|
| | Mean | SD | Mean | SD |
| **MSE** | 382273116.6848 | 3597535.7906 | 395283100.3094 | 1624851.2171 |
| **MAE** | 11109.6831 | 83.1689 | 11412.3168 | 34.0108 |
| **RMSE** | 19551.5897 | 91.9514 | 19881.6858 | 40.8785 |
| **MMRE** | 0.9577 | 0.0072 | 0.9838 | 0.0242 |

# 6.Kemerer Dataset

Kermer consist of 15 samples each sample has 7 attributes.

## feature selection



Pearson Correlation Heatmap

The selected features are KSLOC,AdjFP,RAWFP which are highly correlated features whose correlation coefficient is greater than 0.2



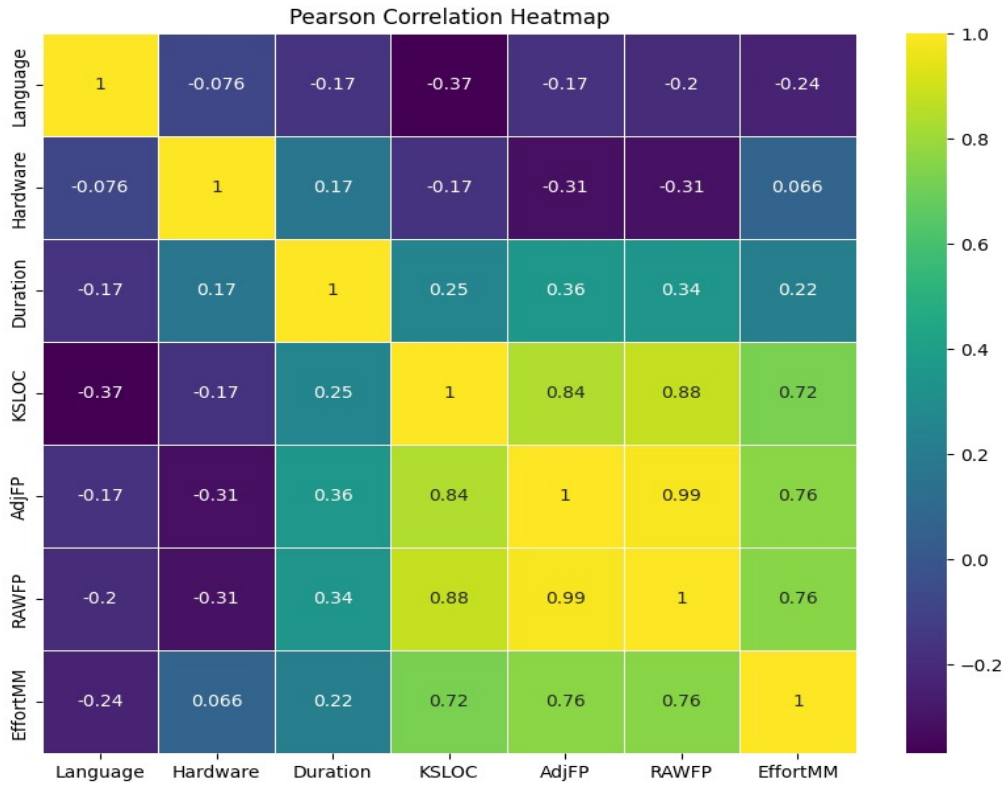Comparison of Mean and Std Dev For Kemerer Dataset With and out Feature Selction

| Measures | Without Feature Selection | | With Feature Selection | |
|----------|---------------------------|------|------------------------|------|
| | Mean | SD | Mean | SD |
| MSE | 44398.5961 | 672.1044 | 45375.6067 | 397.3796 |
| MAE | 189.7734 | 1.8009 | 191.9485 | 1.3858 |
| RMSE | 210.7037 | 1.5969 | 213.0135 | 0.9342 |
| MMRE | 0.9654 | 0.0092 | 0.9765 | 0.0070 |