

Machine Learning (STQD 6024) Assignment 2

Muhammad Hilman Bin Rozaini (P121535)

Objective and dataset description

The objective of the project is to find out the best Multi Linear Regression model to predict the Unified Parkinson's Disease Rating Scale (UPDRS) score, which measures the severity of the Parkinson's disease based on various features.

From the dataset, firstly it contains information from 20 individuals diagnosed with Parkinson's disease and 20 healthy individuals. The dataset also contains various sound recordings, which includes sustained vowels, numbers, words, and short sentences, taken from each subject. Then, a group of 26 linear and time-frequency based features is extracted from the aforementioned recordings. It captures different aspects of voice characteristics such as jitter, shimmer, pitch, and other relevant measurements.

Therefore, the primary objective of this dataset is to develop a predictive model that can predict UPDRS scores based on the extracted voice features. This can be done by analyzing the relationship between the voice features against the UPDRS scores. Thus, the aim is to identify relevant patterns and build model that can assist the diagnosis of Parkinson's disease to gain insights into the Parkinson's diseases' severity and management of the diseases' condition.

Data preprocessing

The dataset contains no null values in all of their variables, and the data types of these variables are either numerical or integer. Thus, data preprocessing is not required in this dataset since it is clean.

Next, the data is split into 80:20 ratio of training data and testing data respectively. This is needed to test the best model and test the consistency of the model in training data and testing data.

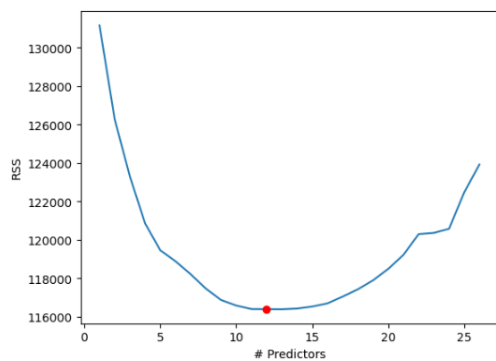
Result summary

Forward selection method

The method selected for predicting the UPDRS score is by using forward selection method as it offers incremental improvement, reduces complexity, provides interpretability, and computationally efficient to select relevant predictors.

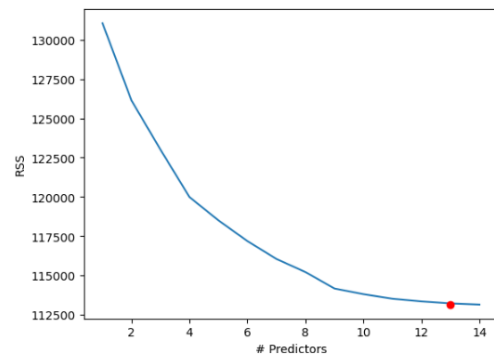
```
plt.plot(models_train["RSS"])
plt.xlabel('# Predictors')
plt.ylabel('RSS')
plt.plot(models_train["RSS"].argmin(), models_train["RSS"].min(), "or")
```

[<matplotlib.lines.Line2D at 0x2a52899adf0>]



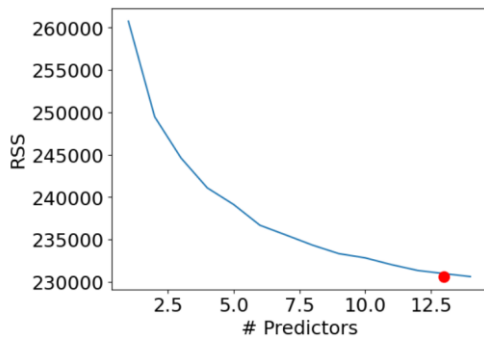
```
plt.plot(models_test["RSS"])
plt.xlabel('# Predictors')
plt.ylabel('RSS')
plt.plot(models_test["RSS"].argmin(), models_test["RSS"].min(), "or")
```

[<matplotlib.lines.Line2D at 0x2a528a12640>]



```
plt.plot(models_all["RSS"])
plt.xlabel('# Predictors')
plt.ylabel('RSS')
plt.plot(models_all["RSS"].argmin(), models_all["RSS"].min(), "or")
```

[<matplotlib.lines.Line2D at 0x15ca0a863a0>]



Based on the dataset, the best MLR model is 13 predictors in both training and testing data. In full dataset as well, the best model consists approximately 13 predictors.

Mean Squared Error (MSE):

```
models_train["RSS"]
```

```
1    131151.778631
2    126268.815874
3    123298.119762
4    120857.707340
5    119462.441125
6    118891.988475
7    118217.167977
8    117469.598595
9    116878.460013
10   116589.791119
11   116409.844934
12   116402.931241
13   116399.275240
14   116435.207129
15   116538.272429
16   116703.035615
17   117056.258499
18   117433.667239
19   117901.744142
20   118495.050732
21   119224.725691
22   120303.843844
23   120367.163649
24   120578.051271
25   122472.958545
26   123922.005604
Name: RSS, dtype: float64
```

```
models_test["RSS"]
```

```
1    131067.471904
2    126141.589957
3    123010.316233
4    119988.248075
5    118502.868451
6    117181.587727
7    116045.172305
8    115210.050312
9    114162.595288
10   113803.111088
11   113514.101632
12   113343.324204
13   113213.329838
14   113138.404551
Name: RSS, dtype: float64
```


relationship between the features (dependent variables) and the UPDRS score (independent variable). Plus, it further indicates that this model is **more robust** and reliable as it does not limit to the specific variations present in the train data, and thus, it can effectively handle the diverse and unseen data better.

Furthermore, the model is able to **avoid overfitting** concern from the RSS argument as mentioned above since it is not overly sensitive to the specific characteristic of the train data. This further indicates that the model's prediction are more accurate and aligned with true values where the model's estimated regression line fits the test data point better, which results in lesser variability in prediction to reduce prediction errors.

Adjusted R-squared:

Adjusted R-squared takes into account the number of predictors in the model, and penalise the inclusion of unnecessary features. Since the adjusted R-squared is higher in test data as well, it suggest that the model is able to explain a larger proportion of the variance in the test data while utilising a more **parsimonious predictors**. Thus, it indicates that the model has effectively identifies and included the most informative predictors which results to a more concise model.

Similar to R-squared, a higher value in test data suggest that the model can **generalise well** beyond the train data while managing the complexity of the predictors. It is **more robust** to handle the diverse and unseen data while maintaining the model from over-fitting which affect the interpretability.

Higher adjusted R-squared in test data is crucial as it improve the model prediction accuracy, and that it is realiable for prediction of new and unseen data. Thus, this further suggest that the model's estimated regression line fits the test data point better which results to reducing biasness and improve prediction accuracy.

In summary, a higher adjusted R-squared value in the test data suggests that the model is more parsimonious, generalizes well beyond the train data, improves prediction accuracy, and demonstrates robustness and general applicability. This highlights the model's ability to provide meaningful insights and reliable predictions when applied to new, unseen data.

Akaike Information Criterion (AIC)

A lower AIC value indicates that the model's performance has a better balance between model fit and complexity. However, the AIC for test data is higher than train data, which suggest that the model's performance in terms of balancing good fit and model complexity is relatively poor for test data.

It indicates that the train data may overfit. It occurs when the model captures noise or random fluctuations in the training data that leads to poor generalisation of unseen data in test data. This also suggests that the model might be too complex or may have included unnecessary predictors that do not contribute to a better fit on new data.

When there is a lack in generalisation, it indicates that the model may not be able to capture the relationships between the predictors against the UPDRS score as it relies on specific characteristics or pattern present in the training data.

Conclusion

As a whole, the model has a positive outlook on its MSE, R-squared and adjusted R-squared. And all of these performance indicators suggest that the model is not overfit and has great generalisation to manage the complexity of the predictors. It also suggests that the model is able to be robust in handling the diverse and unseen data while maintaining the model from over-fitting. This will further improve the prediction accuracy as it is not sensitive to a specific characteristics to any of the features in the dataset.

However, although there is a higher AIC value in the model, which concerns that the model may be unstable, and may face with decrease in predictive performance, it should only be noted on this potential issue as the other three performance indicators are proving otherwise observation.

Therefore, a slight precaution is to be noted when applying this model for predicting the UPDRS score.