

# Linear Regression on the BIOSCAD data set

## *Introduction*

### **Overview:**

The aim of the project was to create a model that, based on biomarkers measured on the skin, would predict the SCORAD of a patient. We used the data from the AMC BIOSCAD study.

## *Method*

### **Data sets:**

The AMC BIOSCAD study data is first pre-processed to make it suitable for training – see pre-processing report ([../Preprocessing/preprocessing-report.pdf](#)) for details. The pre-processed data (see pre-processing document) is split in to three variants. All three variants are based on the combined non-lesional data for both AD and control patients. This gave a total of 100 data points. Data from lesional skin was not used because, after pre-processing, there were only 17 data points remaining.

- In the first variant, the continuous variables are logged before normalizing.
- In the second, the continuous variables are not logged before normalizing.
- The third variant only uses a subset of the first variant. Since in the original data, many values are below detection range or below fit curve range, we only use variables within the detection range: IL-1a, IL1 $\beta$ , age and the categorical variables (FLG, skin type, sex).

### **Models:**

For data sets 1 and 2, elastic net regularization with varying values of both alpha and lambda was used. The data sets were split up to use 60% for training, 20% for validation (hyperparameters (alpha and lambda) selection), and 20% for testing.

We then report the model trained with all available data and the weighted mean of the hyperparameters across cross-validation iterations (the weights corresponds to the relative performance of each model).

For data set 3, a simple linear regression was used as the number of variables was greatly reduced. The data set was split in to 80% training and 20% testing. Once again, 100 different splits of the data set were used to see how model performance varied.

Each model outlined above was created twice, once using objective SCORAD as the dependent variable, and once using SCORAD.

#### **Performance evaluation:**

To calculate performance of a given model the root mean square error (RMSE) between the predicted result and the actual result was calculated. We also calculated the accuracy, defined as the frequency of 'successful predictions'. A prediction was considered succesful when the absolute prediction error was less than 9 points for objective SCORAD and 10 points for total SCORAD. These values were based on the minimum clinically important difference (1).

In order to see how the performance could vary depending on how the testing and training data was split, we recorded the performance of each cross-validation iteration and plotted these in a box plot.

To see if the models were better than an average predictor, we computed the p-values for a one-sided Wilcoxon rank-sum test where the alternative hypothesis was that the model accuracy was higher than the average predictor accuracy are shown below. We also computed the p-values for a one sided t- test where the alternative hypothesis was that the model RMSE was lower than the average predictor RMSE are shown below

#### **Feature selection:**

To see which biomarkers are important, we plotted the coefficients of the final, optimised, model using a bar graph. As all of the attributes are normalised before training, these coefficients can be used to directly compare the importance of each attribute to the output SCORAD.

#### **Residual analysis:**

We also created a plot of the residuals after training and testing on the whole data set to allow us to check the distribution assumptions of our linear regression.

#### **Average predictor:**

We also created a model that simply predicted the average value of either total or objective SCORAD. This was used as a null model. To compare the models, we computed both the kappa coefficient and the p-value of a one sided Wilcoxon rank-sum test where the alternative hypothesis was that the model accuracy was higher than the average predictor accuracy.

## Results

### SCORAD:

	Logged (Elastic net)	Unlogged (Elastic net)	Reduced (GLM)
Alpha	0.376	0.624	-
Lambda	1.62	2.34	-
RMSE	13.65	27.60	14.74
RMSE (p-value)	$1.91 * 10^{-20}$	1	$3.85 * 10^{-2}$
Accuracy	52.7%	49.0%	48.6%
Accuracy (p-value)	$1.77 * 10^{-21}$	$5.49 * 10^{-7}$	$5.47 * 10^{-4}$
Kappa coefficient	0.263	0.206	0.017
Average predictor (RMSE/Accuracy)	25.9	56.3	30.3

Table 1

### oSCORAD:

	Logged (Elastic net)	Unlogged (Elastic net)	Reduced (GLM)
Alpha	0.393	0.606	-
Lambda	1.73	2.12	-
RMSE	10.78	15.75	10.75
RMSE (p-value)	$2.23 * 10^{-7}$	1	$4.44 * 10^{-3}$
Accuracy	68.8%	64.9%	60.8%
Accuracy (p-value)	$8.77 * 10^{-2}$	$3.88 * 10^{-2}$	$2.34 * 10^{-2}$
Kappa coefficient	0.243	0.148	0.063
Average predictor (RMSE/Accuracy)	15.7	24.3	17.7

Table 2

### Average predictors:

The average prediction models (null models) had the following performance:

	RMSE	Successful predictions
Total SCORAD	15.06	35.8%
Objective SCORAD	10.76	58.8%

Table 3

oSCORAD logged data results:

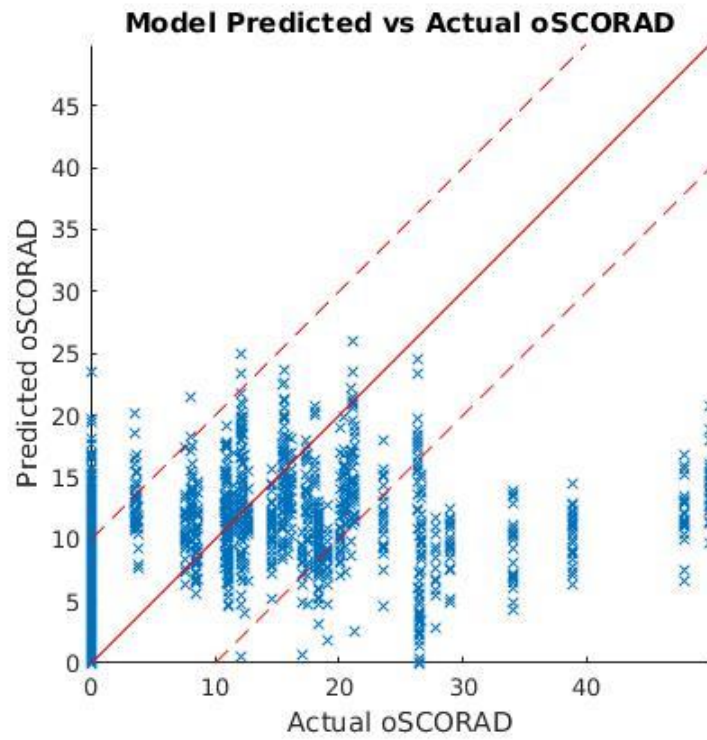


Figure 1

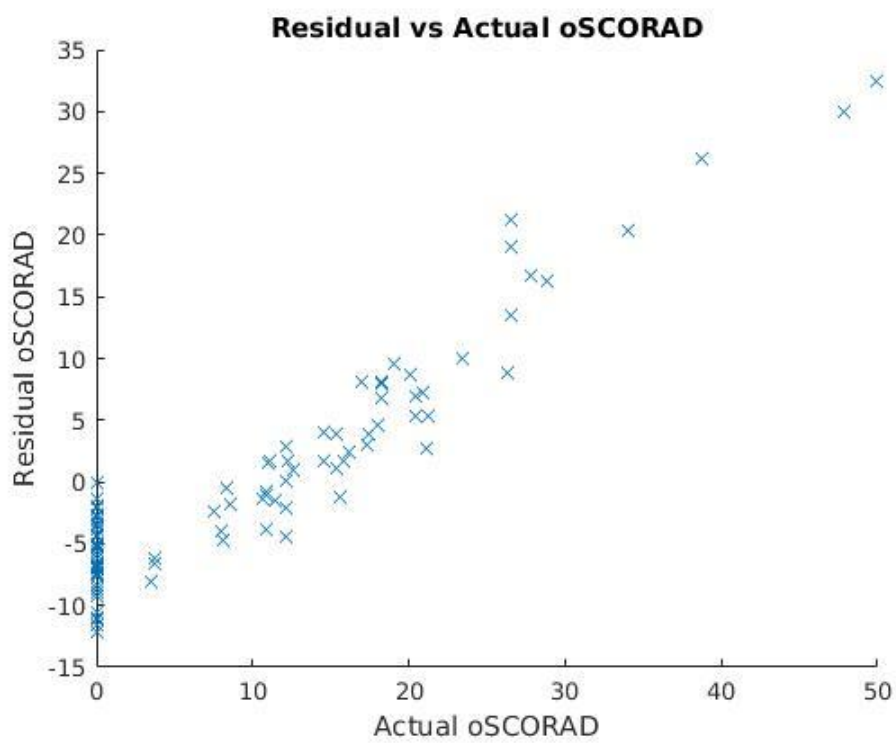


Figure 2

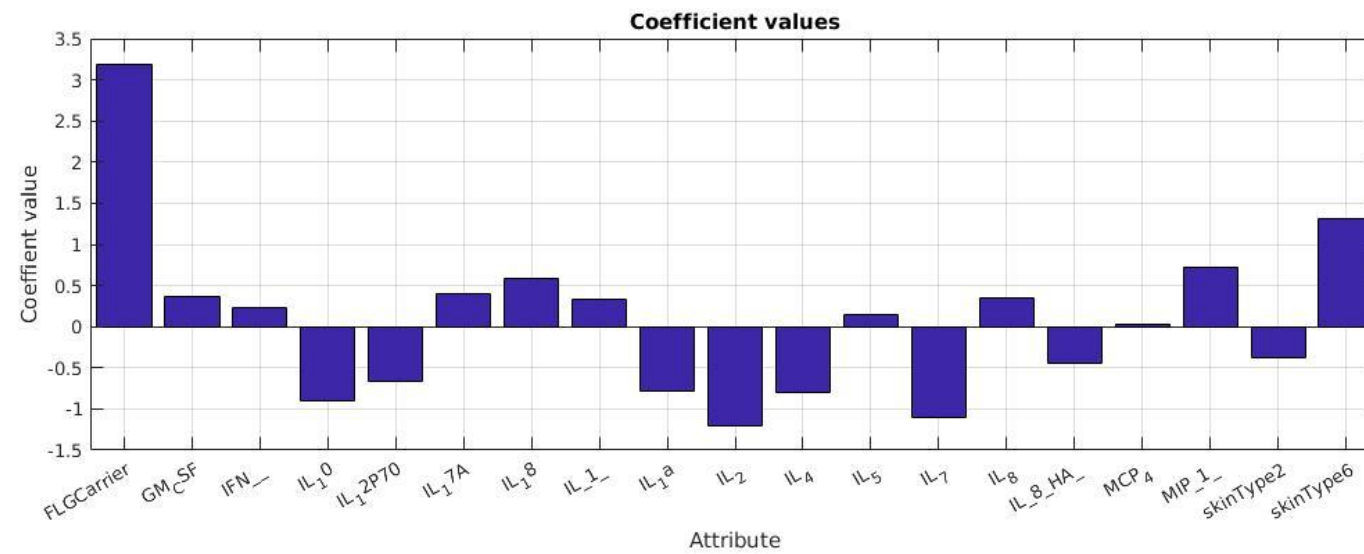


Figure 3

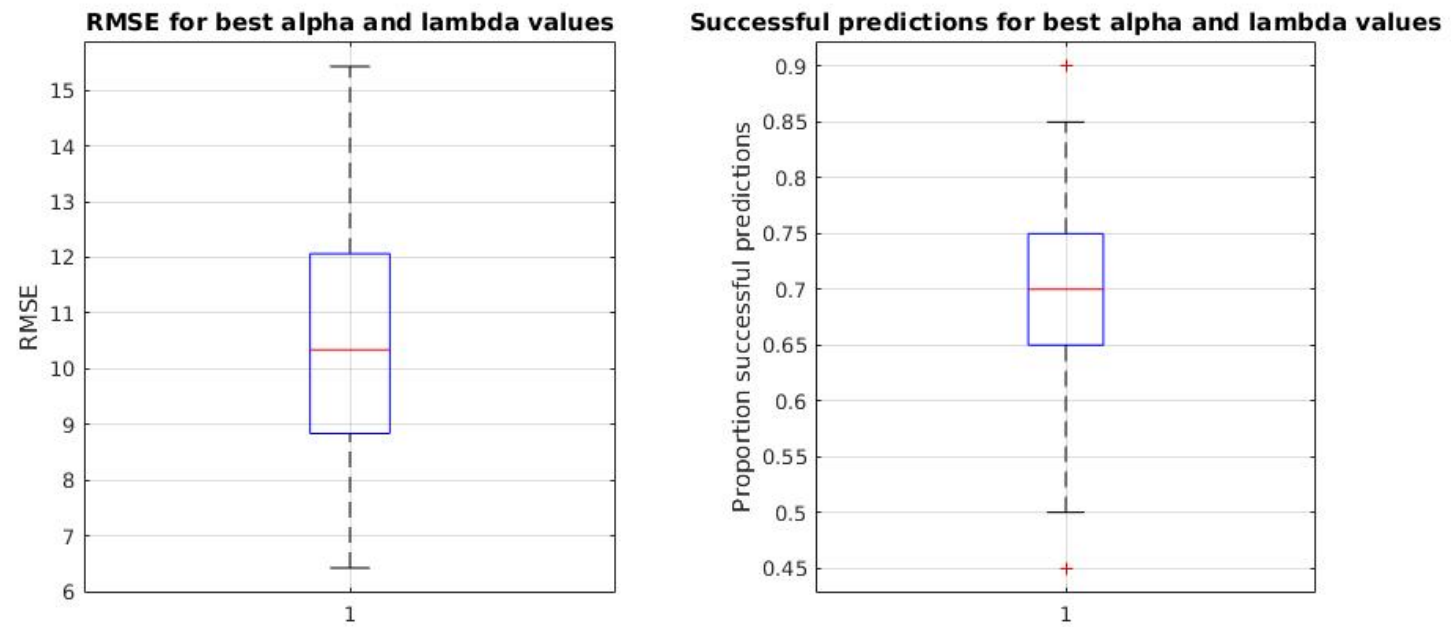


Figure 4

totSCORAD logged data results:

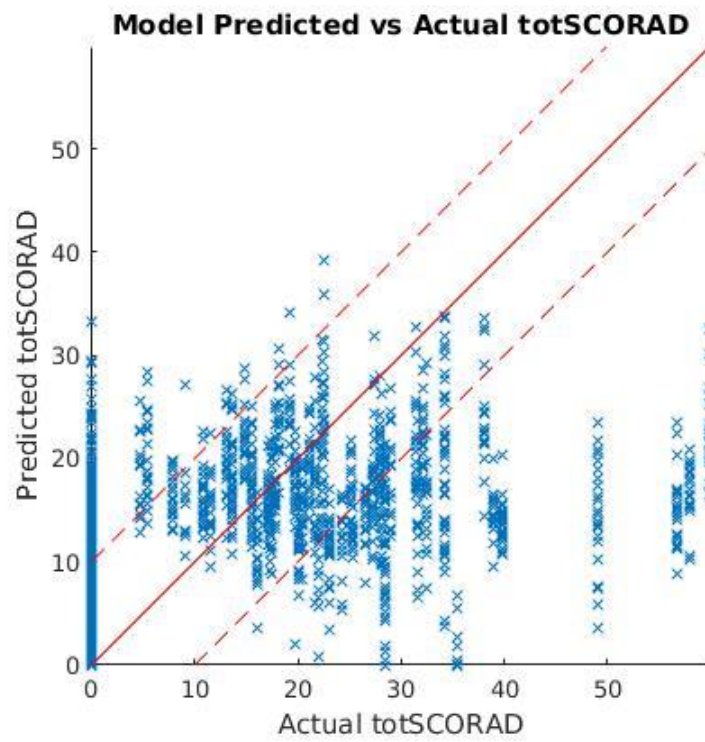


Figure 5

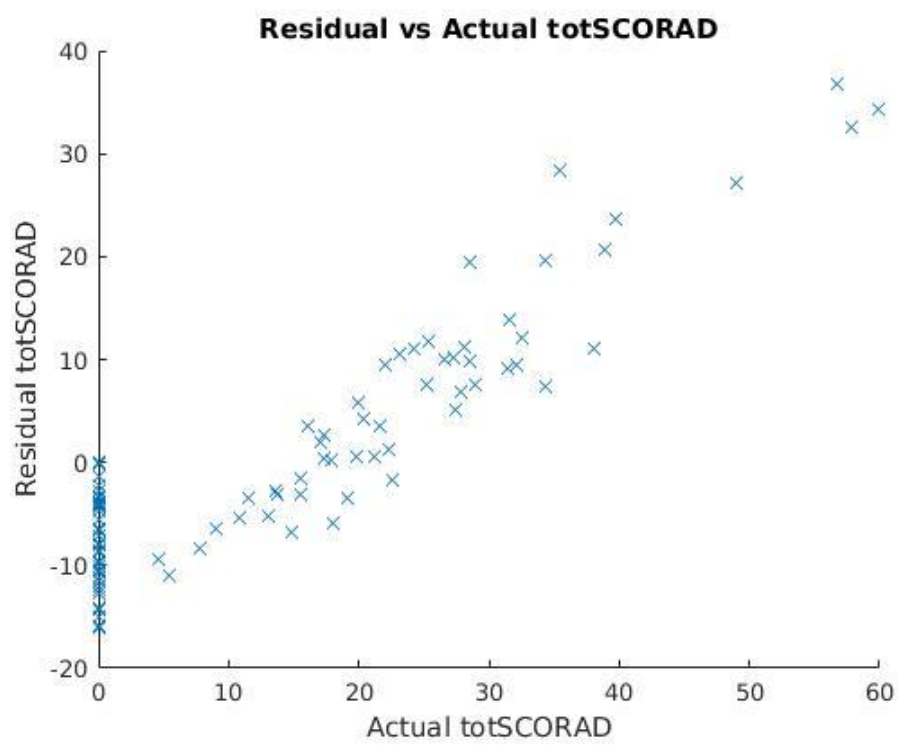


Figure 6

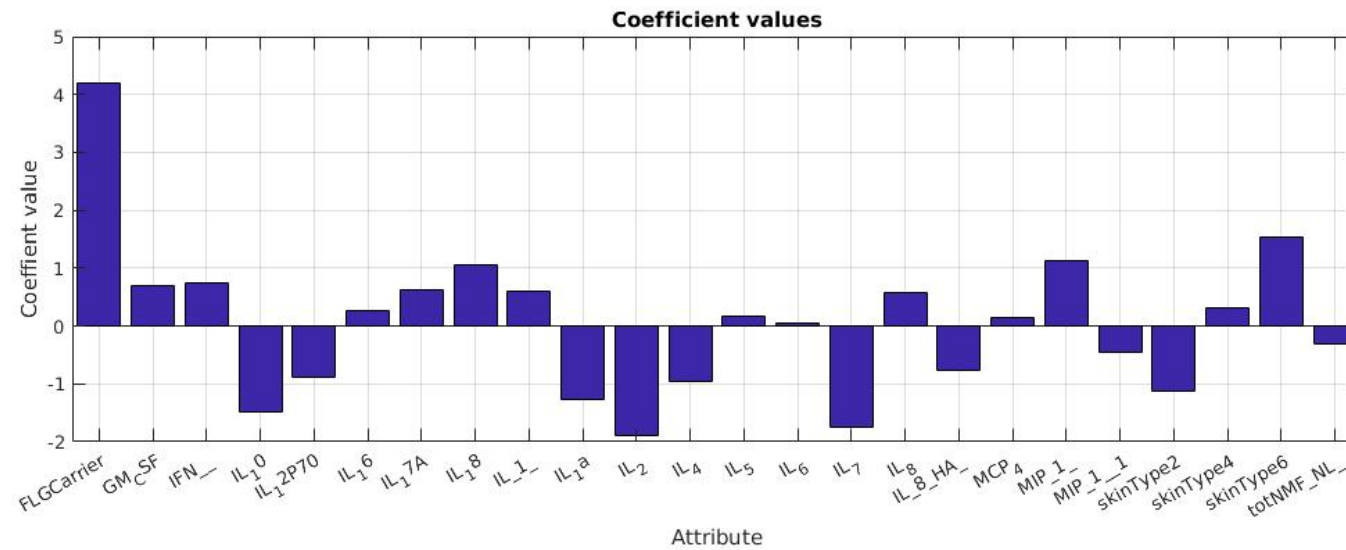


Figure 7

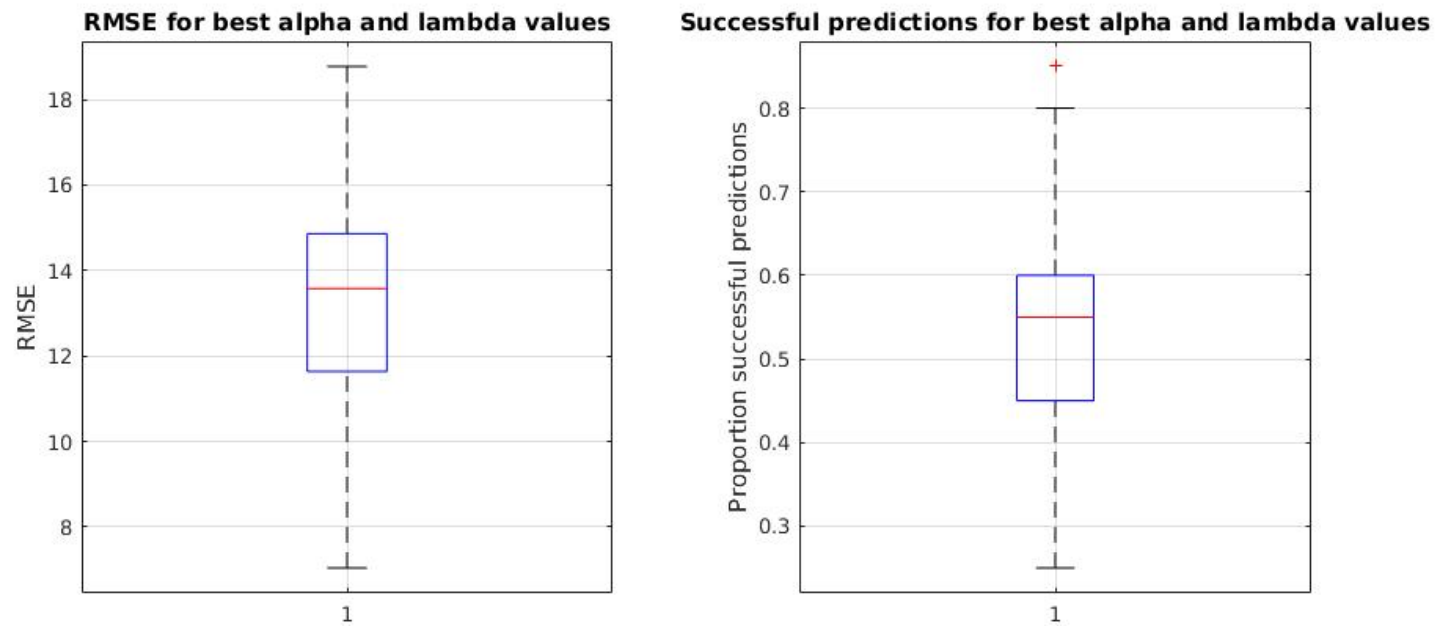


Figure 8

oSCORAD unlogged data results:

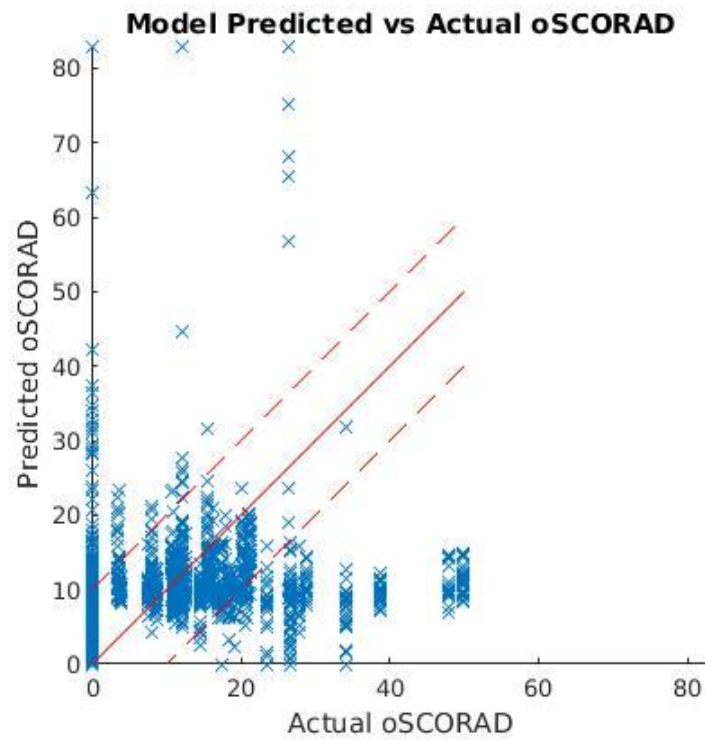


Figure 9

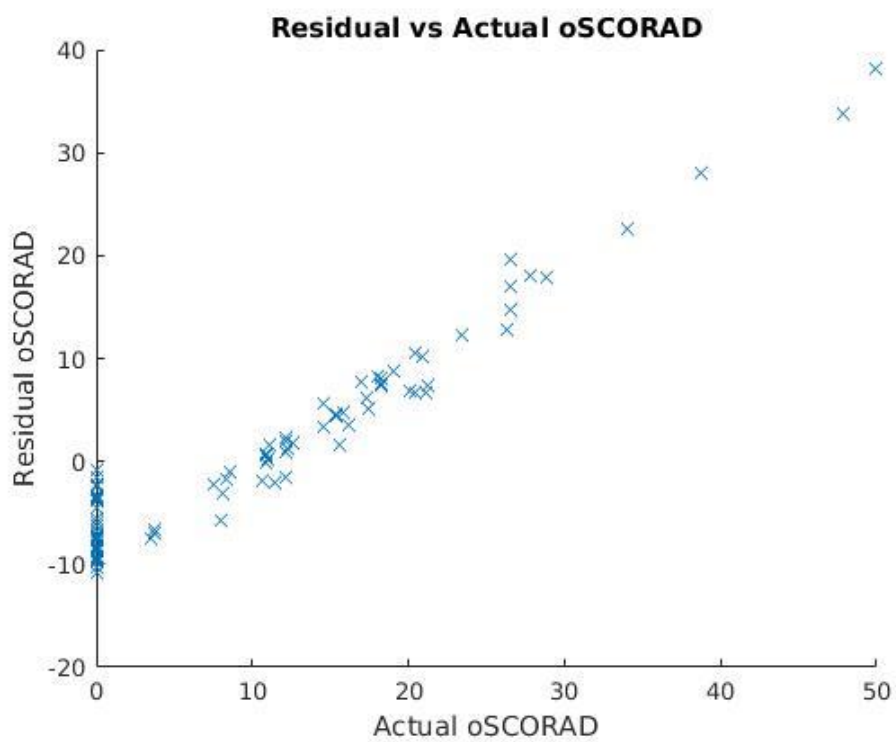


Figure 10



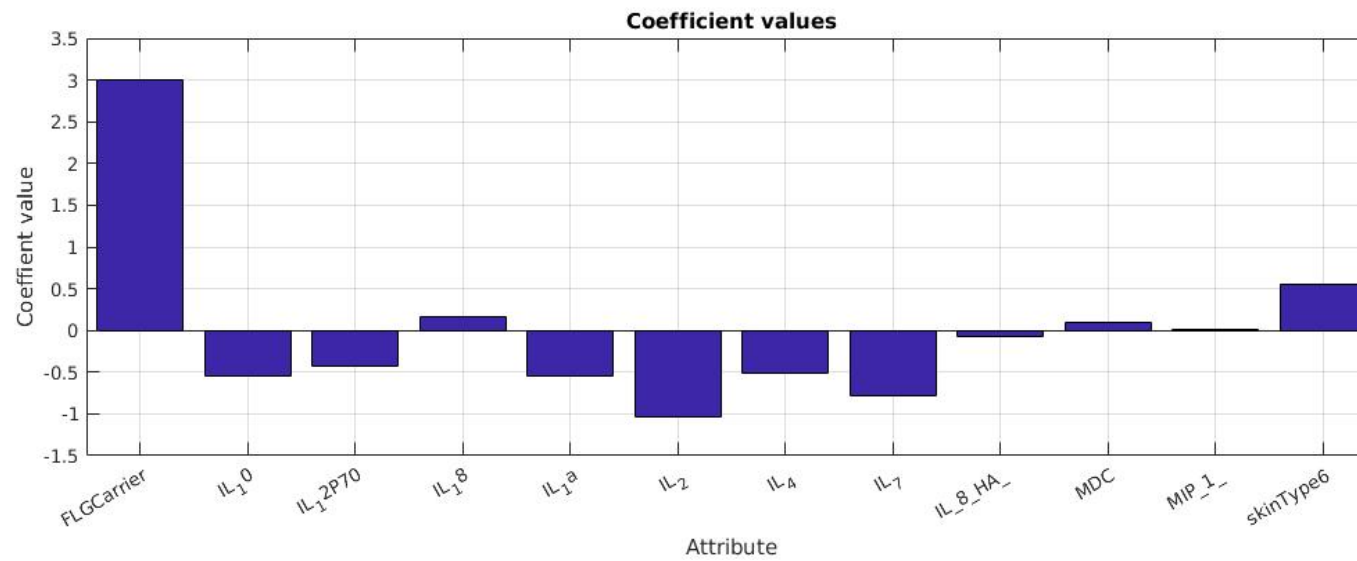


Figure 11

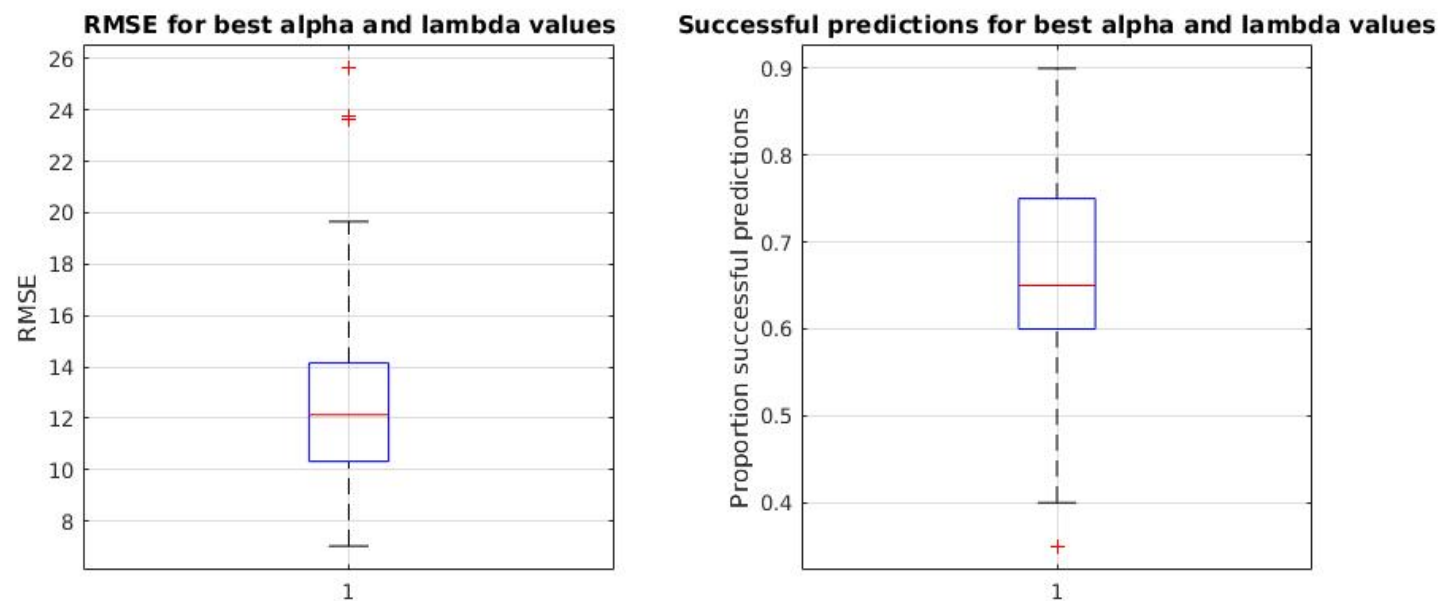


Figure 12

totSCORAD unlogged data results:

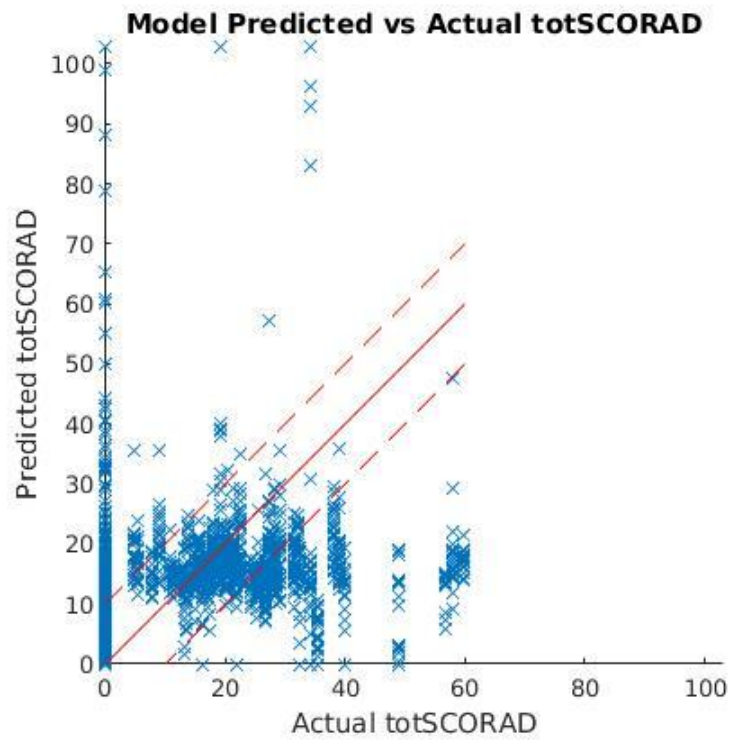


Figure 13

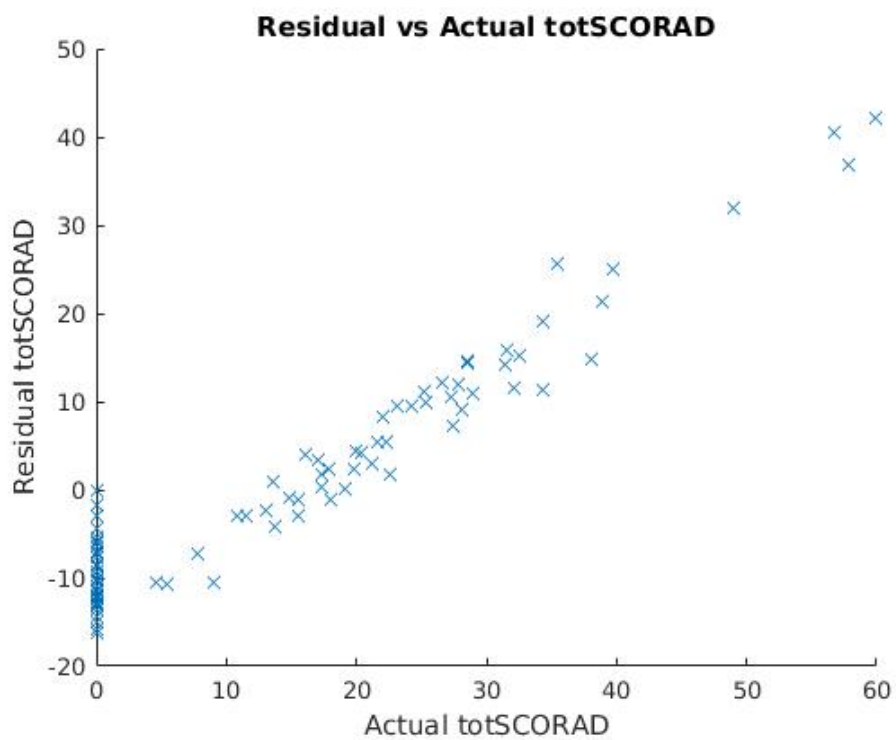


Figure 14

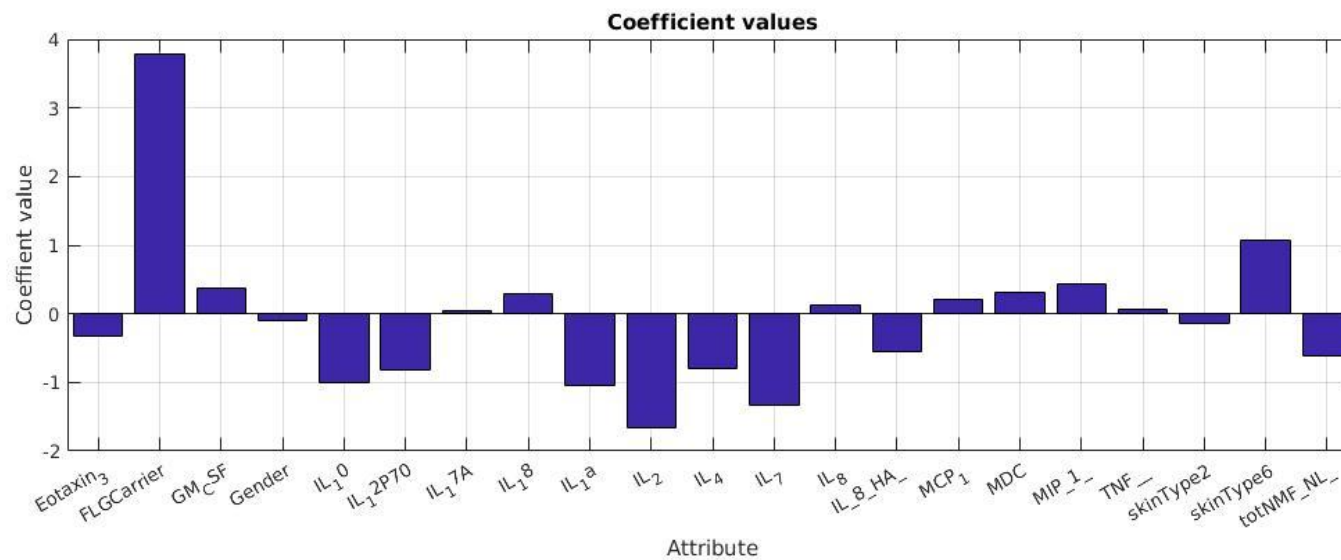


Figure 15

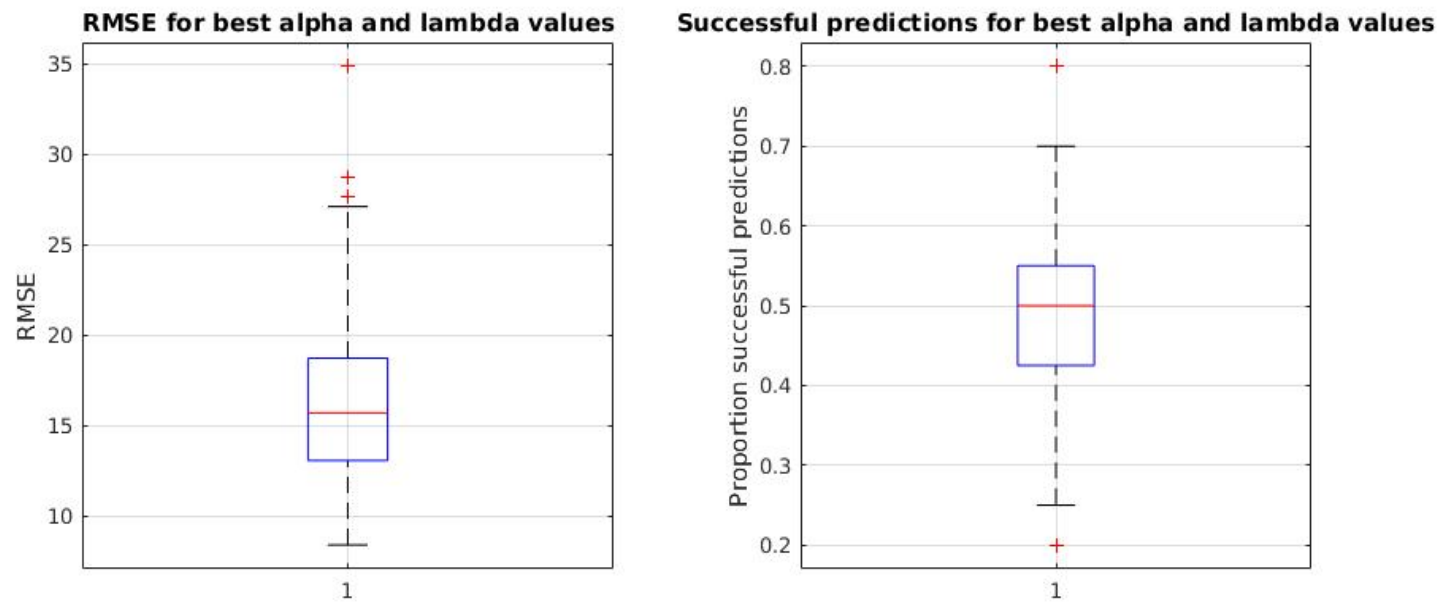


Figure 16

oSCORAD reduced subset data results:

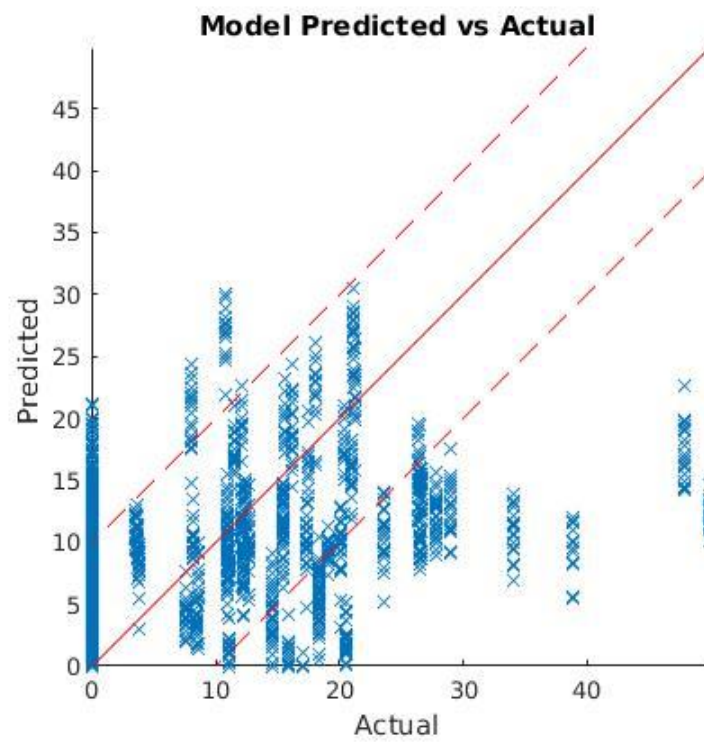


Figure 17

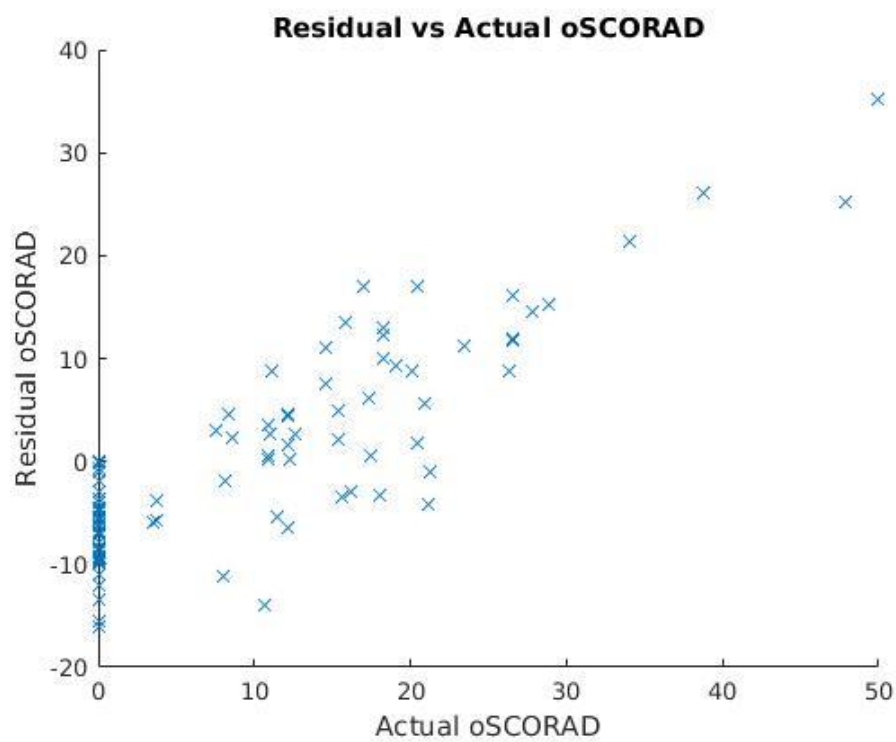


Figure 18

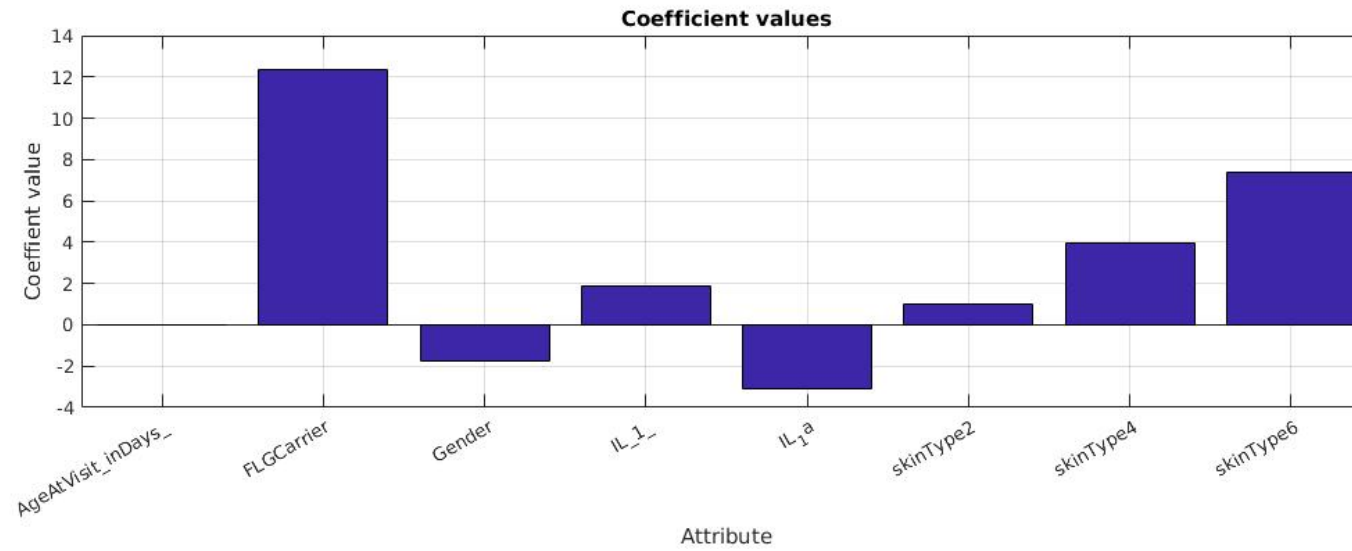


Figure 19

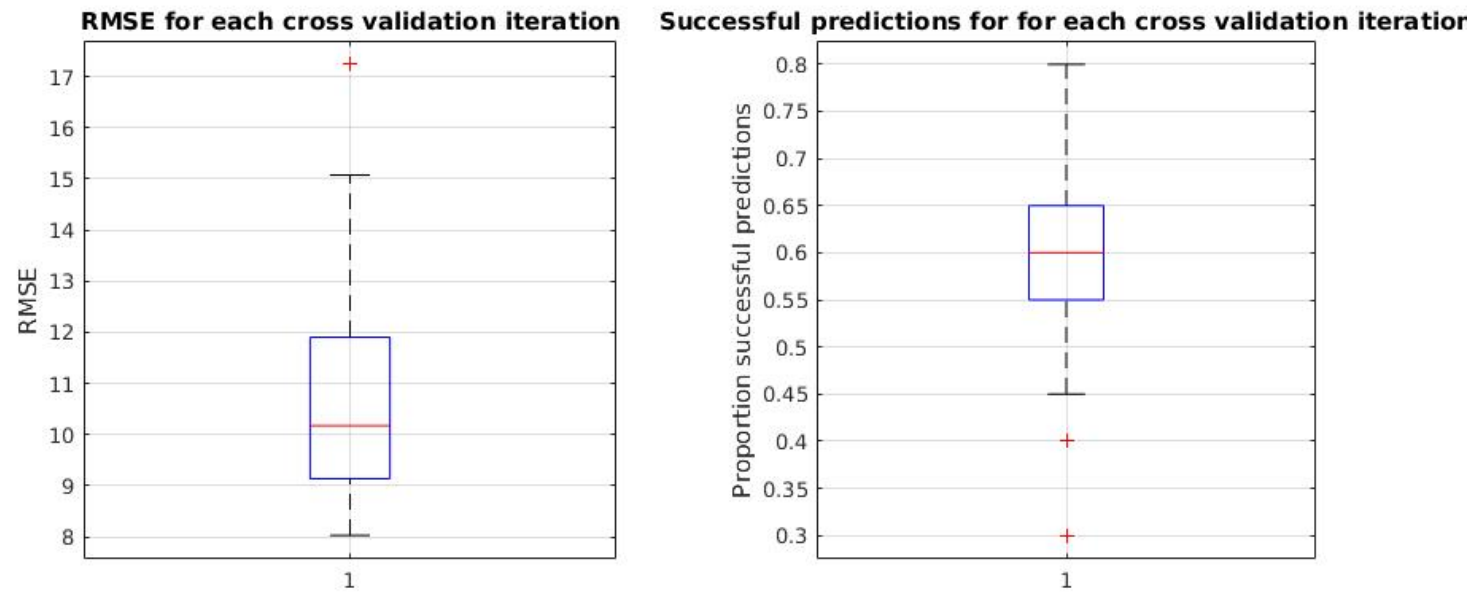


Figure 20

totSCORAD reduced subset data results:

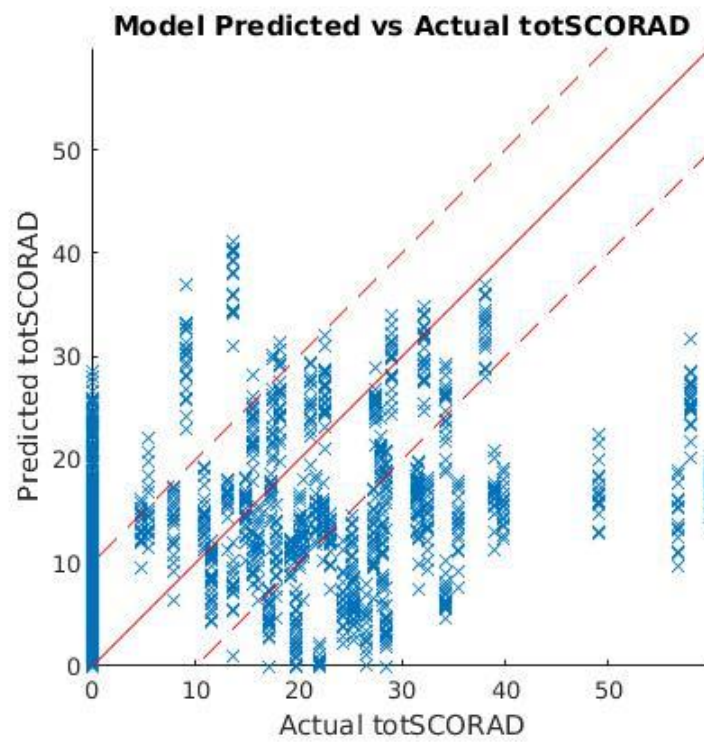


Figure 21

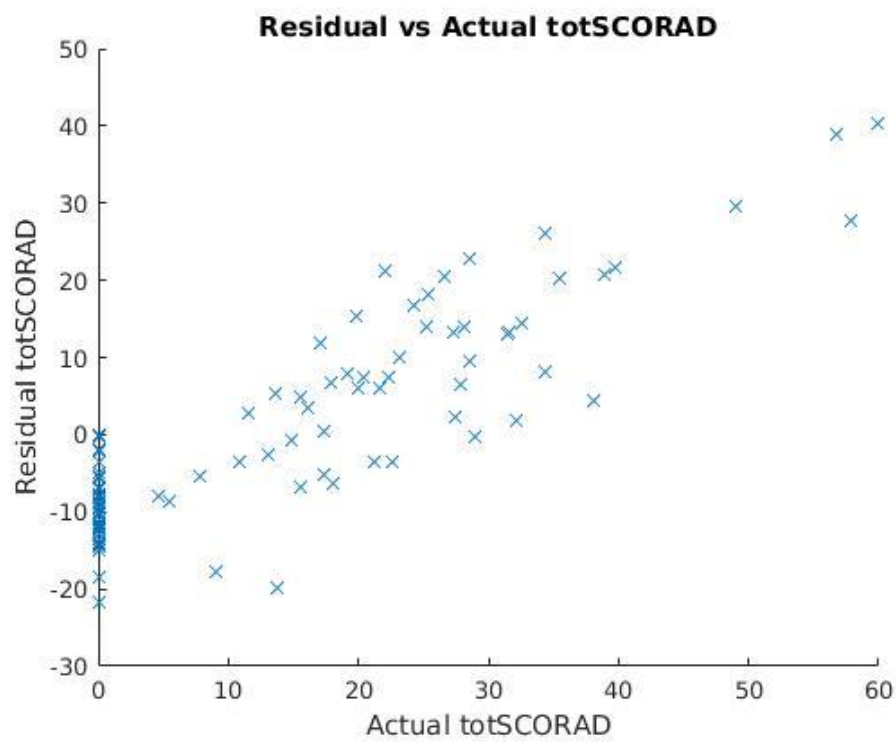


Figure 22

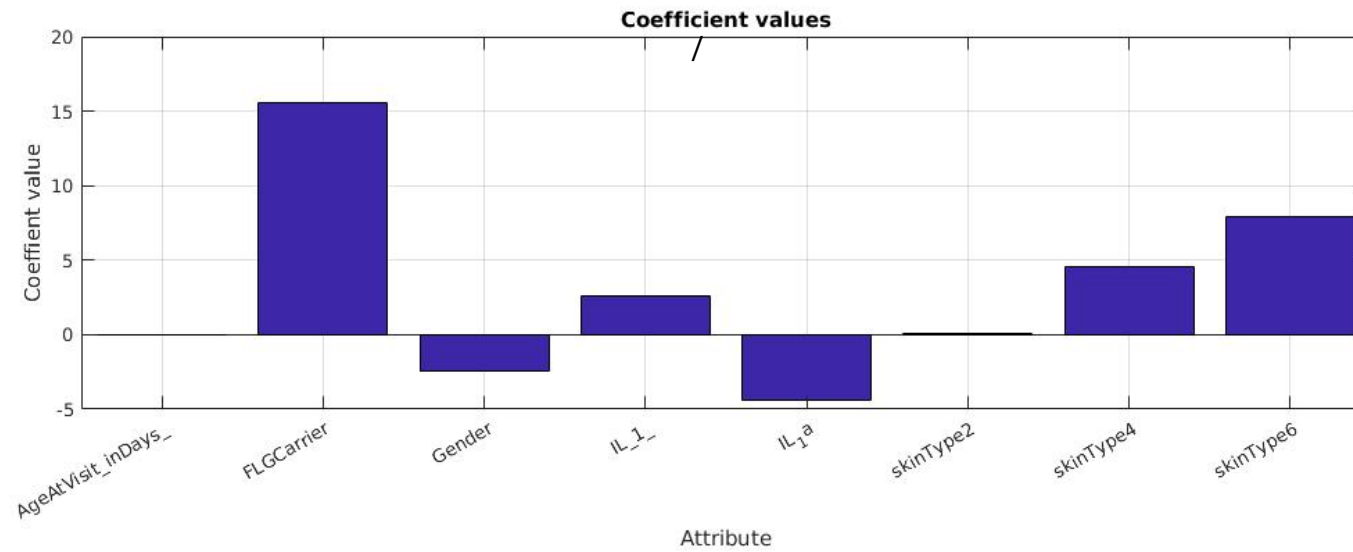


Figure 23

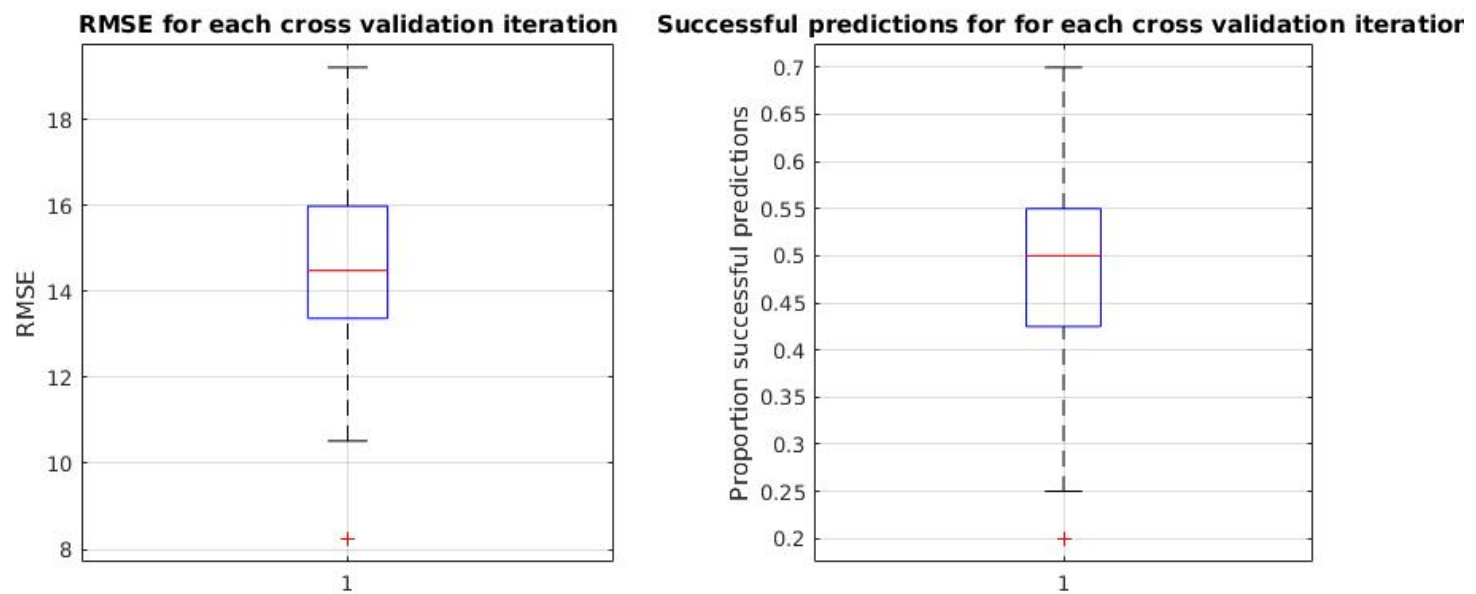


Figure 24

## *Discussion*

### **Performance evaluation:**

As shown above, the models trained on logged data performed better than those based on unlogged data. It can be because logging the data reduces the effect that large value outliers have. This is particularly important as the training data is small, thereby making it harder to spot outliers. Evidence of this can be seen in the very large predictions made in the unlogged models (figures 9 & 13) – likely a result of outliers in the training data.

In addition, models using objective SCORAD as the dependent variable performed better than those using SCORAD.

### **Coefficients:**

From the coefficient graphs, we see that the most important input was found to be the presence of an FLG mutation (FLG carrier input) which, when present, increased the SCORAD predictions by approximately 3.2.

Other biomarkers which had the most significant positive coefficients are skinType6 (1.3), MIP\_1 (0.7), and IL\_18 (0.59). Biomarkers which had the most significant negative coefficients are IL\_2 (-1.2), IL\_7 (-1.1), and IL\_10 (-0.9).

Most coefficients have very little impact on SCORAD, as a result, the prediction range is quite small.

### **Residuals:**

The residuals for each model shows a linear increase in residual value with SCORAD. This indicates that the input data alone cannot accurately predict SCORAD and that there are missing attributes that have a large impact on SCORAD.

### **Comparison to an average predictor:**

As seen from the p-values in tables 1 & 2, some of the models perform significantly better than the average prediction model. However, we see that models that perform better with RMSE do not necessarily perform better with percentage successful predictions and vice-versa. This is probably due to small numbers of points which had extreme differences between the predicted and actual values affecting the RMSE more than the percentage successful predictions.

Overall, this means that although the other models may perform better than the null model, as they still have low accuracy values themselves (tables 1 & 2), they are still not good models.



The lack of performance is to be expected as the coefficients are small so the prediction stays within a small range regardless of input.

This is because, as mentioned earlier, in all but two of the continuous attributes the majority of the data was below the detection range and, therefore, not accurate. As the model is only as good as the input data (garbage in garbage out principle), it is not surprising to have poor performance.

In addition, the data set was too small. Of the 100 data points, only 53 were from patients with atopic dermatitis. Once testing and validation sets were extracted, only 60 points remained to train the model with 35 attributes.

## *Bibliography*

### **Bibliography**

1. Bos JD, Schram ME, Spuls IP, Leeflang MM, Lindeboom R, Schmitt J. EASI, (objective) SCORAD and POEM for atopic eczema: responsiveness and minimal clinically important difference. *European Journal of Allergy and Clinical Immunology*. 2011 September; 67(1).