

Homework 3
Applied predictive modelling
Supriya Bachal

6.1. Infrared (IR) spectroscopy technology is used to determine the chemical Makeup of a substance. The theory of IR spectroscopy holds that unique Molecular structures absorb IR frequencies differently. In practice a spectrometer Fires a series of IR frequencies into a sample material, and the device Measures the absorbance of the sample at each individual frequency. This Series of measurements creates a spectrum profile which can then be used to Determine the chemical makeup of the sample material.

A Tecator Infratec Food and Feed Analyzer instrument was used to analyze 215 samples of meat across 100 frequencies. A sample of these frequency profiles is displayed in Fig. 6.20. In addition to an IR profile, analytical chemistry Determined the percent content of water, fat, and protein for each sample. If we can establish a predictive relationship between IR spectrum and fat Content, then food scientists could predict a sample's fat content with IR Instead of using analytical chemistry. This would provide costs savings, since Analytical chemistry is a more expensive, time-consuming process:

(a) Start R and use these commands to load the data:

> library(caret)

> data(tecator)

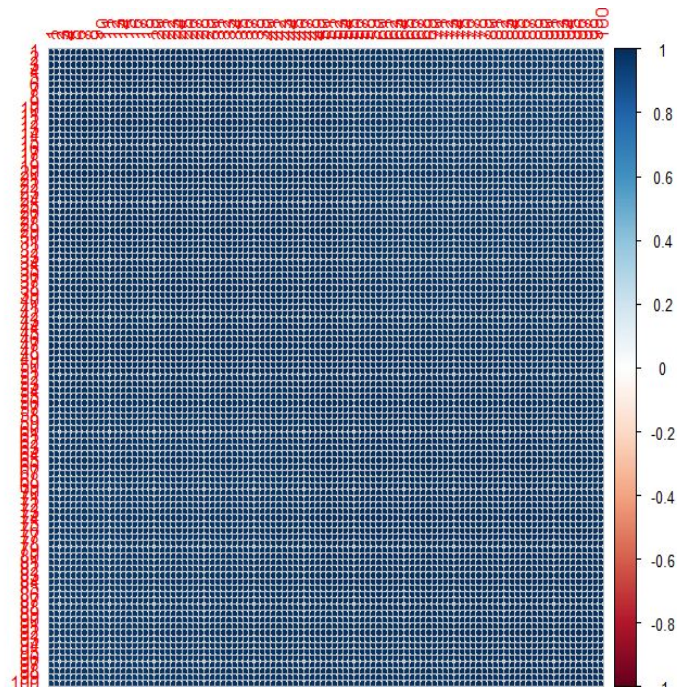
> # use ?tecator to see more details

The matrix absorb contains the 100 absorbance values for the 215 samples,

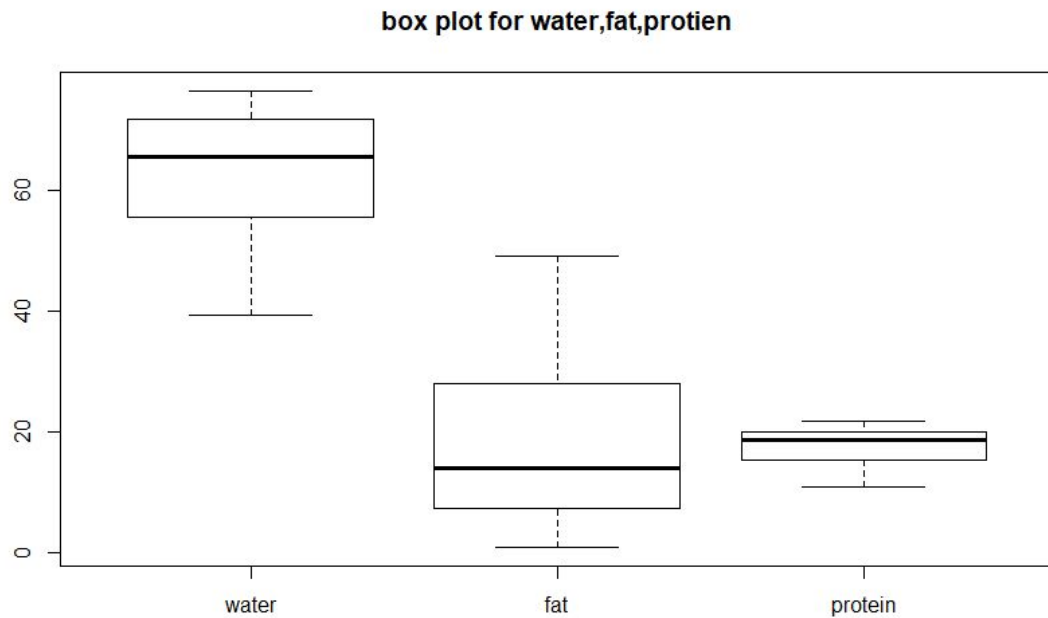
While matrix endpoints contains the percent of moisture, fat, and protein

In columns 1–3, respectively

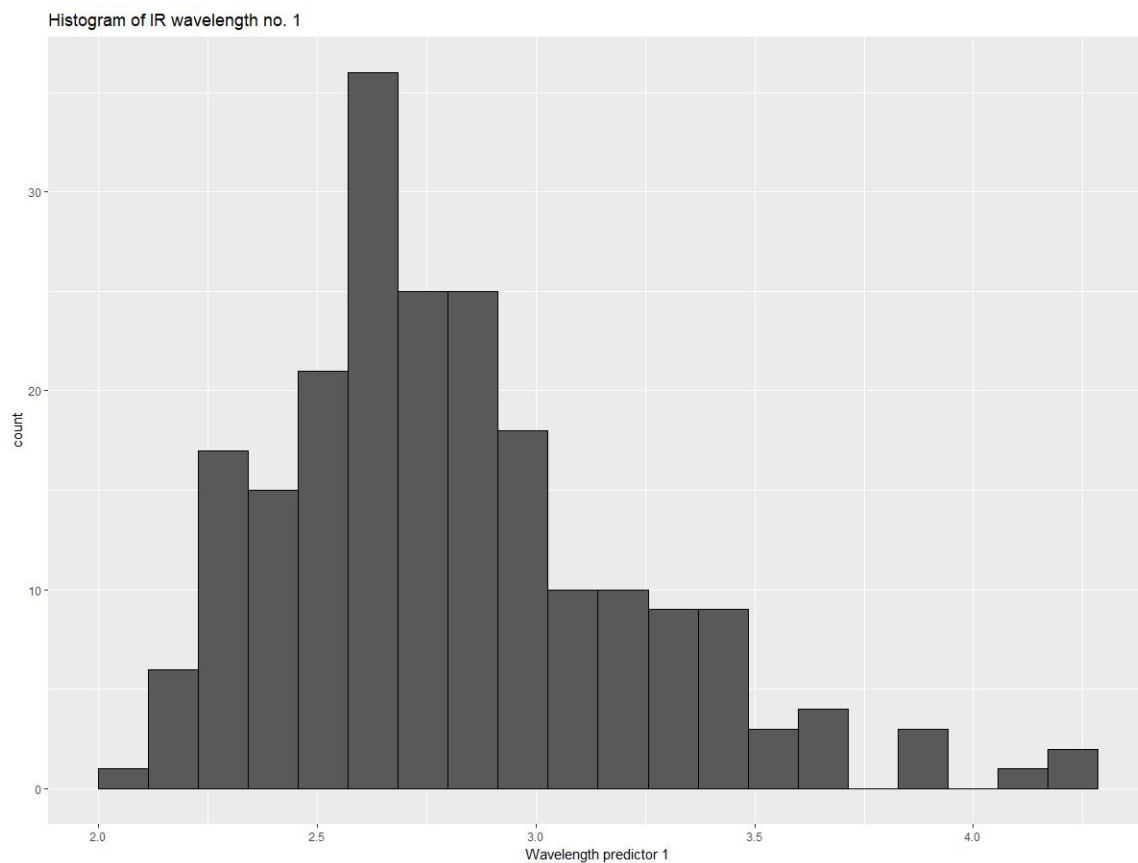
As indicated by the theory of IR, the remarkable structure of atomic retains the distinctive range. Thus, since the fat assimilates particular range given the hypothesis, one can anticipate how much fat is on the nourishment in light of the prescient model, which will spare the time and cost.



The initial Exploration with the dataset confirmed that the all the predictor values in the absorp were very highly correlated with each other. Most Absorption values lie between 2 and 5.



The Box plots for the endpoints leads us to believe that most of the meats have fat contents in the range of 5 to 30. With the median being at around 10. The minimum value for fat is 0.9 and the maximum value of fat is 76.6. We can also interpret the absence of outliers, which means we do not need to use any outlier preprocessing for this data set.



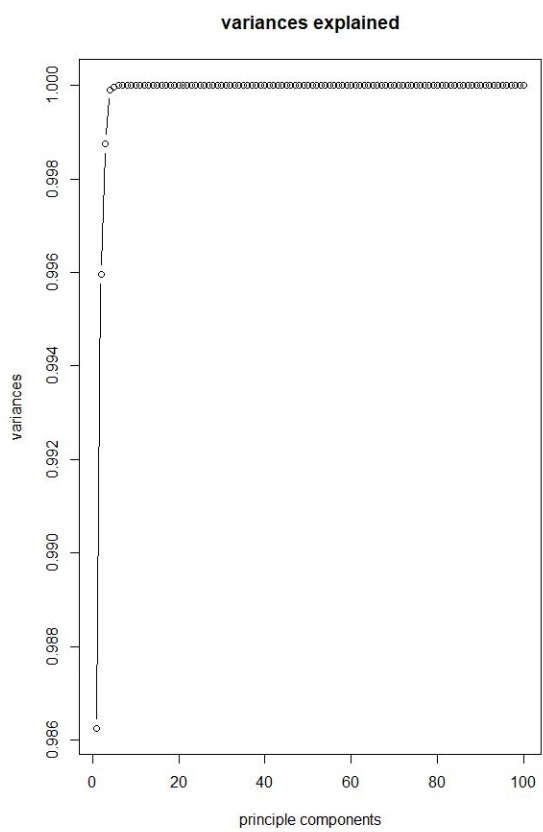
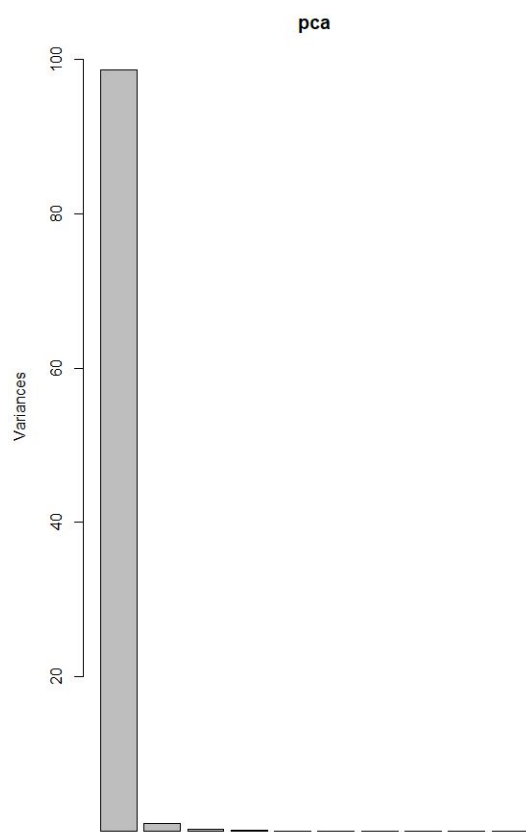
The histogram of the wavelength shows us there is a skewed pattern in the predictors. The predictors are positively skewed; however, since our models will not be heavily affected by this, we do not transform these predictors.

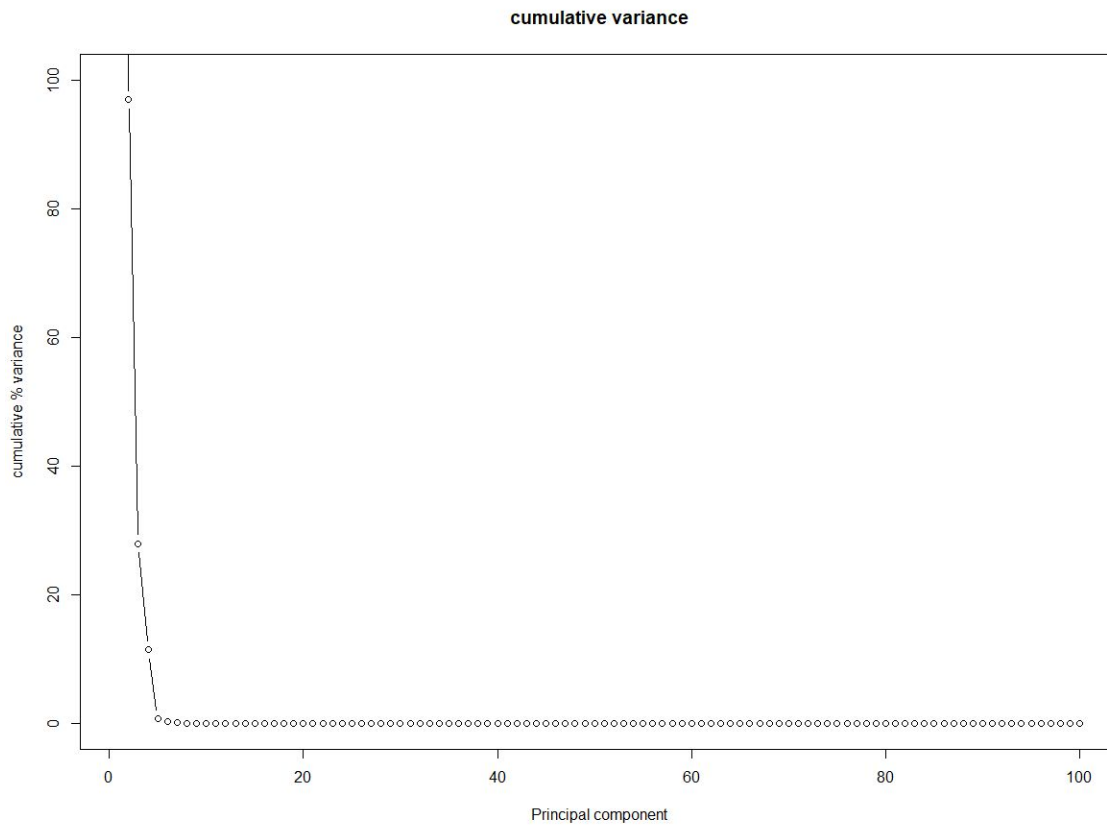
(b) In this example the predictors are the measurements at the individual frequencies.

Because the frequencies lie in a systematic order (850–1,050 nm), the predictors have a high degree of correlation. Hence, the data lie in a smaller dimension than the total number of predictors (100). Use PCA to determine the effective dimension of these data. What is the effective dimension?

To carry out the PCA, we use the `prcomp` function with the `caret` package and center and scale the data.

```
pca <- prcomp(absorp, center = TRUE, scale = TRUE)
```



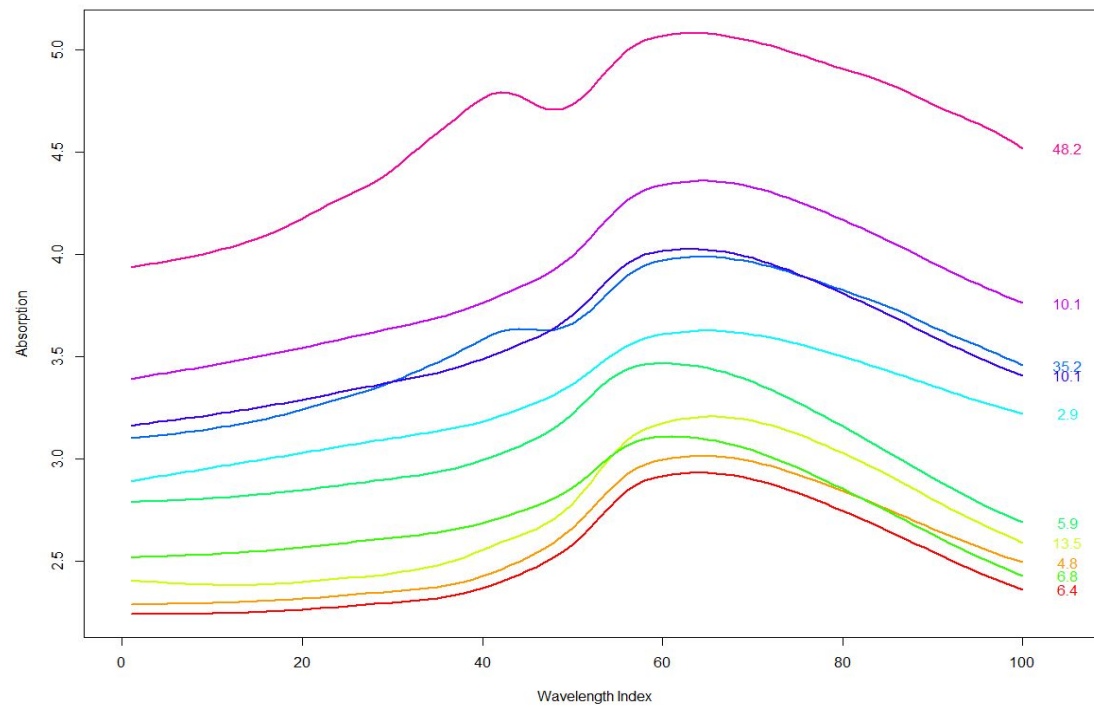


Given the plots, the first principal components explains most of the variance in the data. In addition, after 5 components, the variance explained reached 100%. It means there is high cor linearity among variables. The first plot for visualization was the one with the variance explained however, I could not clearly see the no of components as they were clustered together so I made the second visualization to add clarity. The histogram confirms that after 5 components the variance reaches 100 percent. The summary for the first 5 components is as follows:

	PC1	PC2	PC3	PC4	PC5
Standard deviation	9.931072	0.9847361	0.5285114	0.3382748	0.08037979
Proportion of Variance	0.986260	0.0097000	0.0027900	0.0011400	0.00006000
Cumulative Proportion	0.986260	0.9959600	0.9987500	0.9999000	0.99996000

As seen from the summary 98.62% of the variance is explained by the first component itself and hence it is the major component for the analysis with reference to the linear combination of predictors. Additionally we could also consider nonlinear summarizations of data.

To understand the relationship between the absorption values vs the wavelength, a subset of 10 random values was made and the corresponding, fat values were taken. Every color represents an absorption value and a corresponding fat content between 48.2 and 2.9.



To further investigate the data before splitting it I wanted to see how the fat percentages are distributed within the samples. Here is the table showing the number of instances for the different fat percentages

numbers Freq		
1	0-10	77
2	10-20	61
3	20-30	36
4	30-60	41

c) Split the data into a training and a test set, pre-process the data, and build each variety of models described in this chapter. For those models with tuning parameters, what are the optimal values of the tuning parameter(s)?

The data is not suitable for cross validation since it has only 215 samples. So we decided to use "LGOCV" with 5 repeats. To split the data we use the train control from the caret library.

```
ctrl <- trainControl(method = "LGOCV", number = 5, p = .75)
```

We use 75% as training set and 25% as testing set.

1. The first model built was a Linear Regression model.

To build this model a threshold value of 90% correlation was chosen to select the predictors. We observe from the summary of the model that only one independent predictor was used to build the model. The performance of the model can be evaluated from its RMSE value, which is 11.1.

Linear Regression

163 samples

1 predictor

No pre-processing

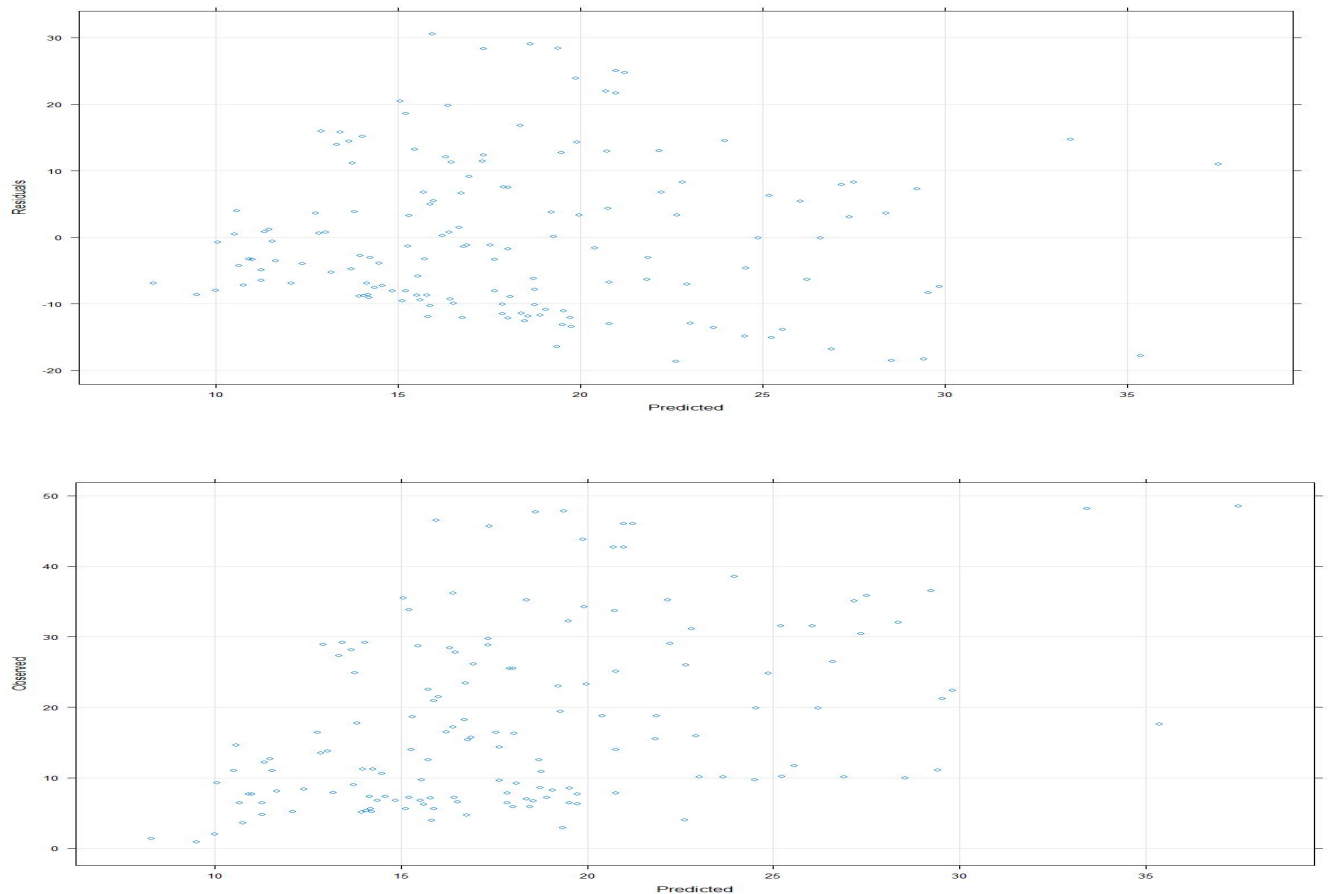
Resampling: Repeated Train/Test Splits Estimated (5 reps, 75%)

Summary of sample sizes: 124, 124, 124, 124, 124

Resampling results:

RMSE	Rsquared	MAE
11.11397	0.2304765	9.14265

These are the two plots to evaluate the performance of the model. One of the plots is the observed value vs the predicted value the other one is that of the residuals. As we can see there is no pattern in the observed vs predicted values which means that there is no or very little correlation between them. Which indicates that the model did not perform very well. As for the residuals, there are more residuals in the bottom portion and higher valued residuals in the top portion however; the residuals are evenly distributed.



2. RLM model is usually used for dealing with outliers and homoscedastic data ,even though our dataset has none of these properties since the question demanded to use all models we implement RLM as well.
For rlm we cannot have a singular predictor covariance matrix thus we preprocess with PCA:

set.seed(1)

```
rlm_model=train(Train[,101],Train[,101],method="rlm",preProcess=c("pca"),trControl
=ctrl )
```

Summary of the RLM model is as follows:

RMSE was used to select the optimal model using the smallest value.

The final values used for the model were intercept = TRUE and psi = psi.hampel.

Robust Linear Model

163 samples

100 predictors

Pre-processing: principal component signal extraction (100), centered (100), scaled (100)

Resampling: Repeated Train/Test Splits Estimated (5 reps, 75%)

Summary of sample sizes: 124, 124, 124, 124, 124

Resampling results across tuning parameters:

intercept	psi	RMSE	Rsquared	MAE
FALSE	psi.huber	21.30030	0.3254849	18.620285
FALSE	psi.hampel	21.31960	0.3178190	18.646044
FALSE	psi.bisquare	21.29260	0.3249492	18.573792
TRUE	psi.huber	10.61796	0.3248411	8.141147
TRUE	psi.hampel	10.55503	0.3200146	8.326236
TRUE	psi.bisquare	10.64477	0.3246944	8.129551

RMSE was used to select the optimal model using the smallest value.

The final values used for the model were intercept = TRUE and psi = psi.hampel.

Which indeed means that the RMSE is 10.5 for this model. It establishes that it performs a little better than the linear model.



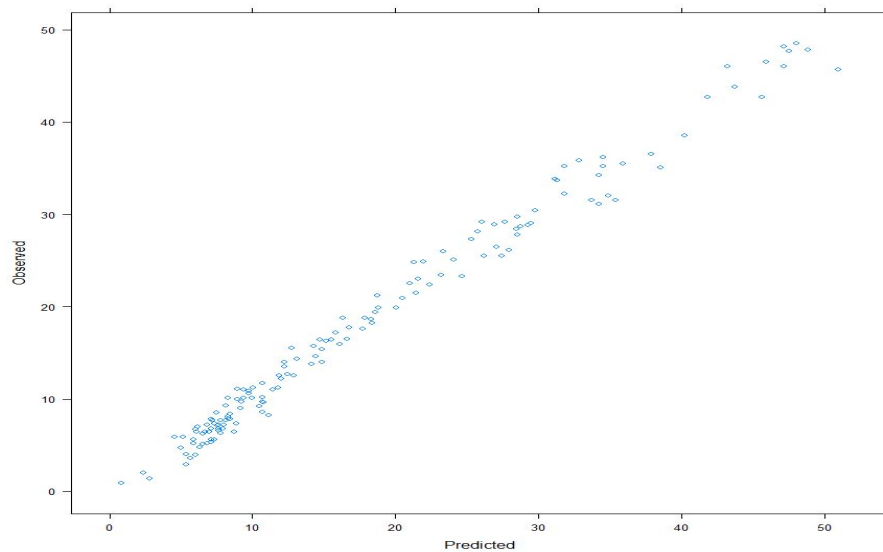
3. The PCR model was implemented on this data using the following code:

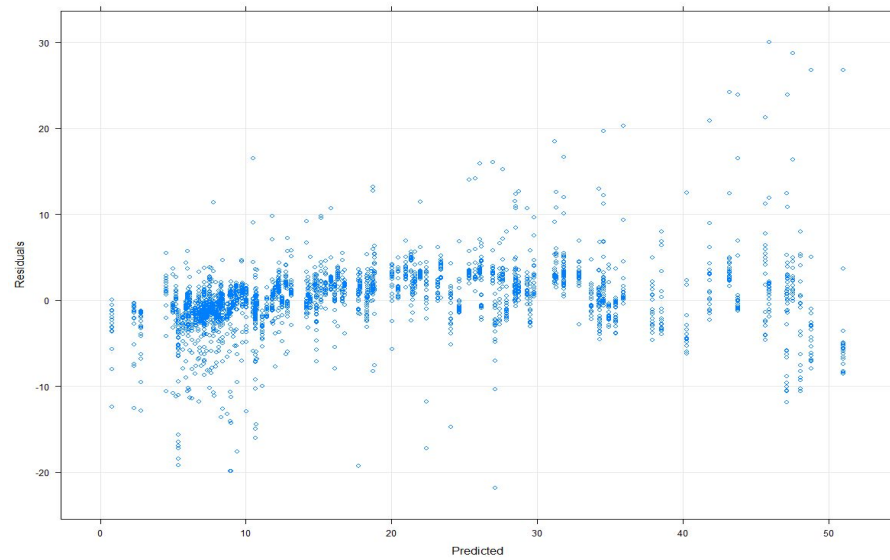
```
PCR <- train(y = Train[,101] , x = Train[,-101], method = "pcr",trControl =  
ctrl,tuneLength=25)
```

The tune length was chosen to be 25 .The best RMSE was 2.58 for 16 predictors. The plots shows that there is a strong pattern in the observed vs predicted values, thus we can say that the model did a good job at predicting the values. The residual plot also shows appropriate behavior.

ncomp	RMSE	Rsquared	MAE
1	10.627729	0.2993532	8.562489
2	10.640628	0.2992878	8.519212
3	8.113044	0.5892546	6.096007
4	4.343164	0.8884487	3.368914
5	3.498749	0.9292436	2.695119
6	3.186437	0.9412674	2.367312
7	3.136030	0.9432869	2.319718
8	3.122363	0.9447836	2.256299
9	3.077401	0.9462370	2.269609

10	3.074187	0.9461419	2.301295
11	2.834915	0.9525326	2.199592
12	2.913620	0.9488626	2.309188
13	2.925539	0.9487389	2.304019
14	2.893150	0.9503614	2.276475
15	2.833587	0.9534902	2.200531
16	2.583810	0.9617547	1.997992
17	2.752988	0.9580129	2.055830
18	2.727945	0.9581410	2.041521
19	2.960649	0.9517838	2.107726
20	2.954541	0.9519300	2.108689
21	2.934897	0.9523863	2.095458
22	2.816743	0.9542456	2.040074
23	2.910588	0.9492736	1.978043
24	2.894319	0.9500870	1.999269
25	2.779262	0.9544376	1.954891



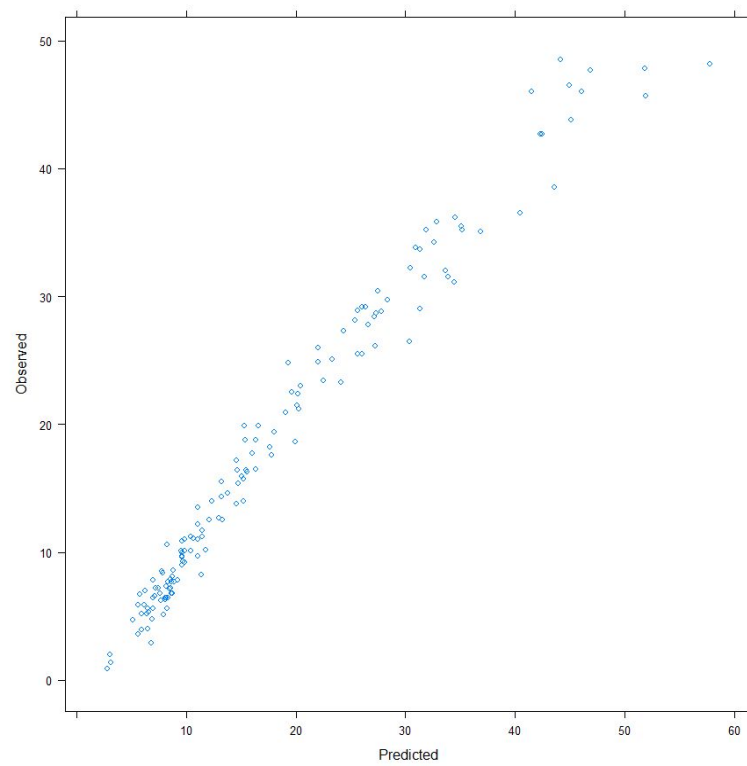
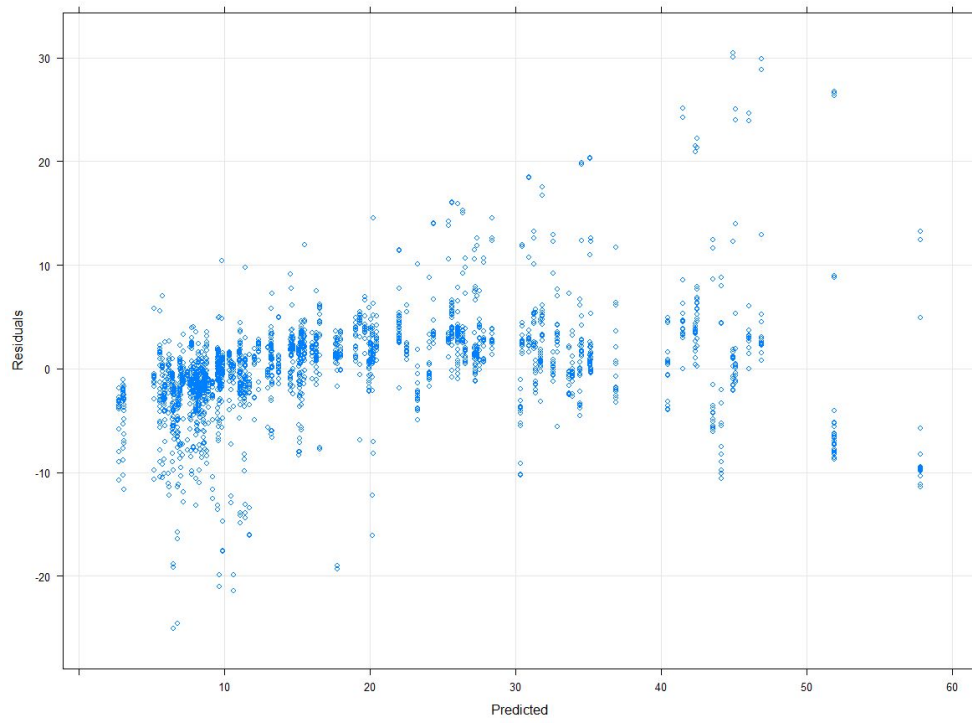


4. PLS Model was implemented using the following code:

```
PLS <- train(y=Train[,101] , x = Train[,-101], method =  
"pls",preProcess=c("center","scale"),trControl = ctrl,tuneLength=25)
```

ncomp	RMSE	Rsquared	MAE
1	10.643722	0.2969975	8.584599
2	6.773286	0.7101878	5.012824
3	5.572418	0.8064884	4.115805
4	4.309526	0.8904063	3.351301
5	3.198210	0.9414655	2.339831
6	3.190998	0.9412995	2.353225
7	3.064207	0.9470725	2.206677
8	3.087945	0.9457460	2.261480
9	2.952283	0.9492857	2.288731
10	2.825851	0.9516280	2.246056
11	2.870393	0.9499581	2.290753
12	2.988820	0.9486453	2.243210
13	2.708144	0.9594643	2.022789
14	2.715399	0.9585977	2.029196
15	2.757921	0.9574485	1.993409
16	2.928037	0.9520755	2.015474
17	2.663415	0.9588269	1.902652
18	2.471738	0.9654171	1.810884
19	2.461300	0.9649078	1.780286
20	2.463416	0.9658904	1.760537
21	2.685104	0.9607576	1.878862

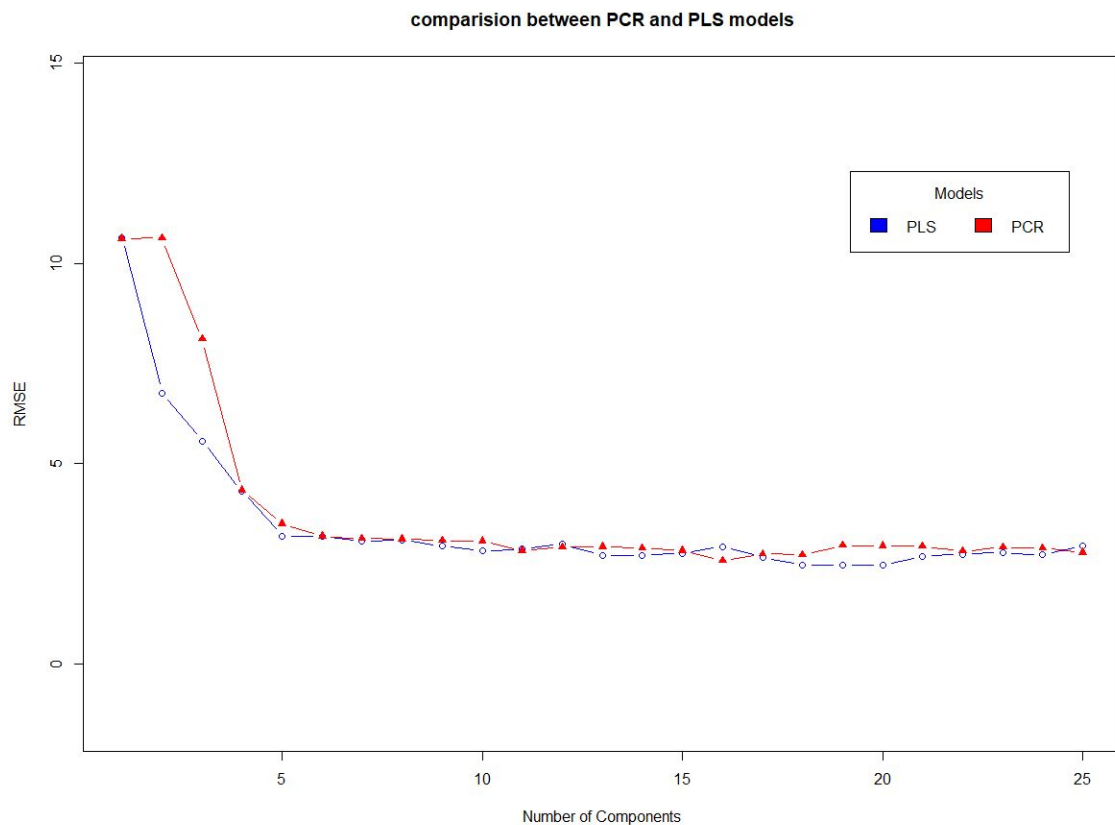
22	2.743685	0.9587711	1.850678
23	2.777816	0.9576040	1.873661
24	2.735002	0.9577879	1.868724
25	2.937455	0.9563290	1.954361



The PLS model was performing similar to the PCR model, as indicated by the plots.

The best value for the RMSE with the PLS model was 2.46 for 19 components.

Next we compare the PCR and PLS model for better performance.



The Graph above shows that the model PCR and PLS perform very similar in terms of RMSE values as 2.58 & 2.46 however the PCR gives a less complex model with 16 predictors as oppose to the 19 predictors use by PLS.

5. Ridge Regression - penalise square of coefficient

Code to implement this model is :

```
RM<- train(y =Train[,101],x = Train[,-101],method = "ridge",trControl =  
ctrl,tuneGrid = ridgeGrid,preProcess=c("center","scale"))
```

The ridge grid was set to :

```
ridgeGrid <- data.frame(.lambda = c(0,0.0001, 0.001, 0.01,0.1))
```

The Ridge regression uses lambda values between 0 to 0.1

Ridge Regression

163 samples

100 predictors

Pre-processing: centered (100), scaled (100)

Resampling: Repeated Train/Test Splits Estimated (5 reps, 75%)

Summary of sample sizes: 124, 124, 124, 124, 124

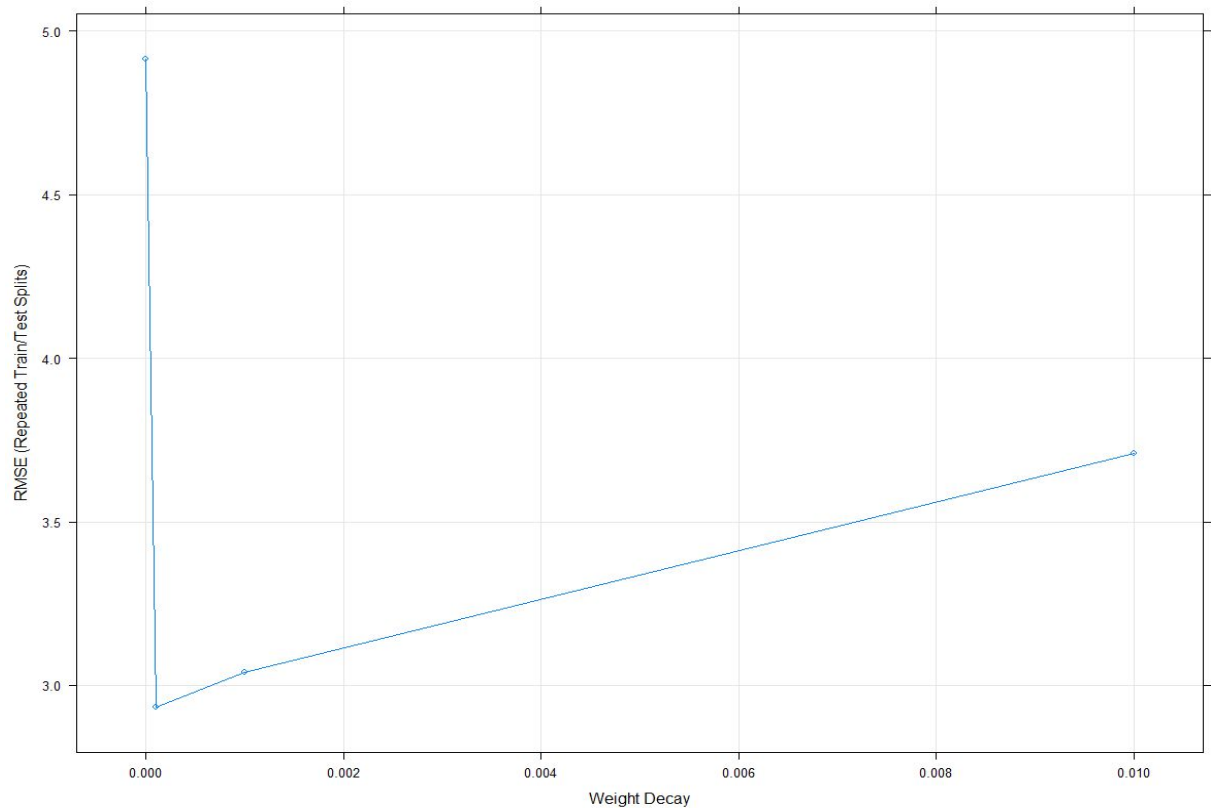
Resampling results across tuning parameters:

lambda	RMSE	Rsquared	MAE
0e+00	4.914019	0.8882616	2.902447
1e-04	2.934167	0.9558180	2.267345
1e-03	3.040376	0.9537380	2.384255
1e-02	3.710175	0.9246771	3.123035
1e-01	4.832154	0.8569034	3.834798

RMSE was used to select the optimal model using the smallest value.

The final value used for the model was lambda = 1e-04.

An interesting observation while implementing the model was that as the tuning parameter is set finer the value of the best lambda moves toward 0. To verify this I implemented a model with Lambda grid extending to 0.00001 and the best lambda value was 0.00001 in this case. There is no static value for lambda it moves closer to 0, as we put in finer parameters.



6. Lasso Model :

The following code was used to implement Lasso model:

```
set.seed(1)
```

```
lasso_grid <- lasso_grid <- expand.grid(.fraction=seq(0,0.1, length=20))
```

```
Lasso <- train(y = Train[,101],x = Train[,-101], method = "lasso", trControl = ctrl,
```

```
tuneGrid = lasso_grid)
```

The grid parameters were selected in the form of fraction parameters between 0 to 0.1.

The lasso

163 samples

100 predictors

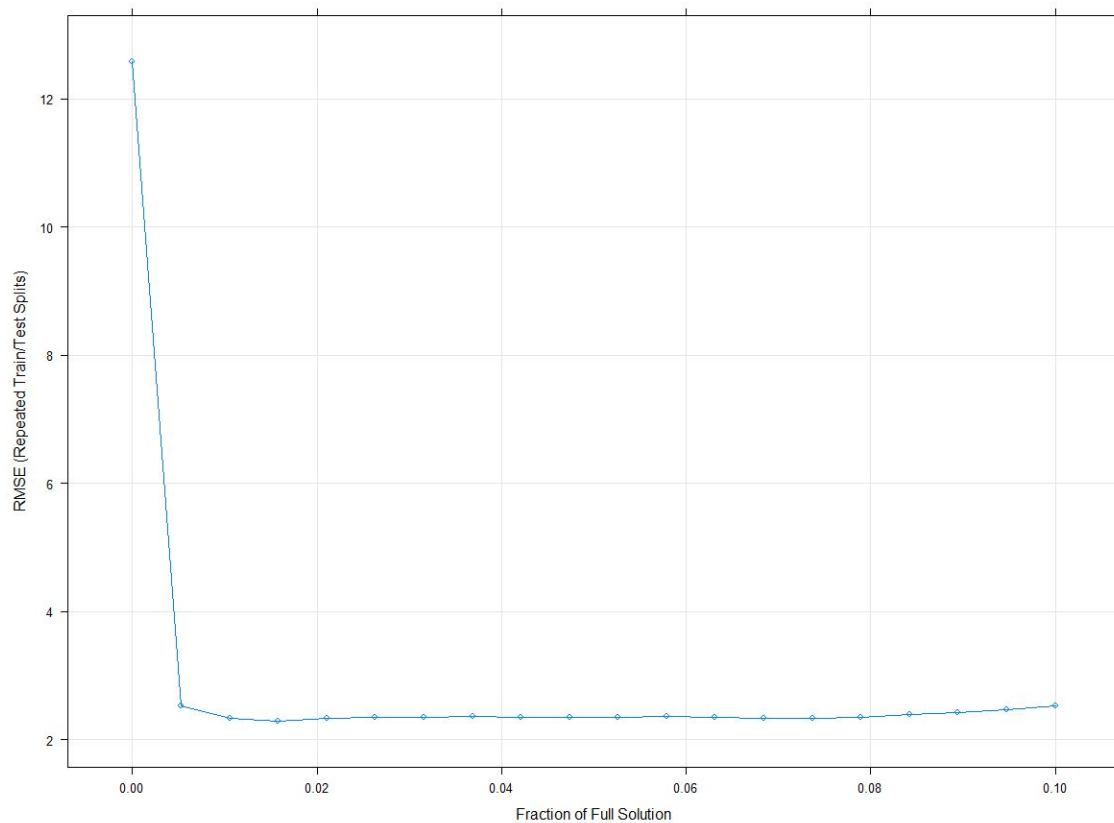
No pre-processing

Resampling: Repeated Train/Test Splits Estimated (5 reps, 75%)

Summary of sample sizes: 124, 124, 124, 124, 124
Resampling results across tuning parameters:

fraction	RMSE	Rsquared	MAE
0.000000000	12.574961	NaN	10.696237
0.005263158	2.526463	0.9622234	1.959691
0.010526316	2.335365	0.9682374	1.794543
0.015789474	2.294162	0.9698386	1.740583
0.021052632	2.329030	0.9691821	1.720583
0.026315789	2.347864	0.9686924	1.691038
0.031578947	2.354168	0.9684607	1.681249
0.036842105	2.358874	0.9684357	1.687682
0.042105263	2.355805	0.9687439	1.694593
0.047368421	2.348086	0.9690578	1.696465
0.052631579	2.347767	0.9691707	1.696021
0.057894737	2.361642	0.9688745	1.698470
0.063157895	2.347330	0.9695324	1.694744
0.068421053	2.331335	0.9701243	1.686742
0.073684211	2.332103	0.9702683	1.681896
0.078947368	2.353445	0.9699895	1.696900
0.084210526	2.391985	0.9692924	1.717262
0.089473684	2.428490	0.9684846	1.731387
0.094736842	2.472501	0.9675176	1.744057
0.100000000	2.523552	0.9662329	1.757509

RMSE was used to select the optimal model using the smallest value.
The final value used for the model was fraction = 0.01578947.



7. Elastinet model (combination of ridge and Lasso):

The code used for implementing the Elastinet model is as follows:

```
enet_grid <- expand.grid(.fraction=seq(0,0.1, length=20), .lambda = c(0,
0.0001, 0.001, 0.01))
Enet <- train(y = Train[,101],x = Train[-101], method = "enet", trControl = ctrl,
tuneGrid = enet_grid)
```

A grid function with fraction between 0.001 to 0.1 and lambda values between 0 to 0.0001 are used.

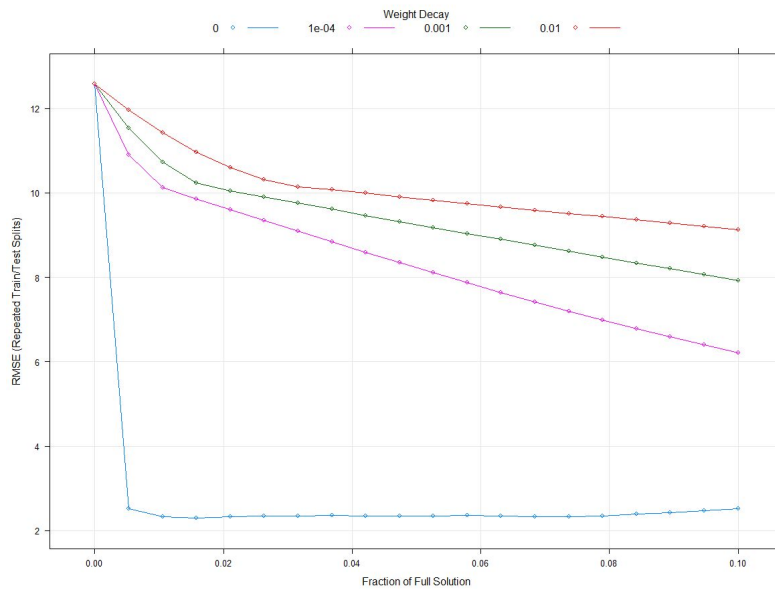
The model generates a grid for variety of combinations of these values(see appendix) but the best value selected are :

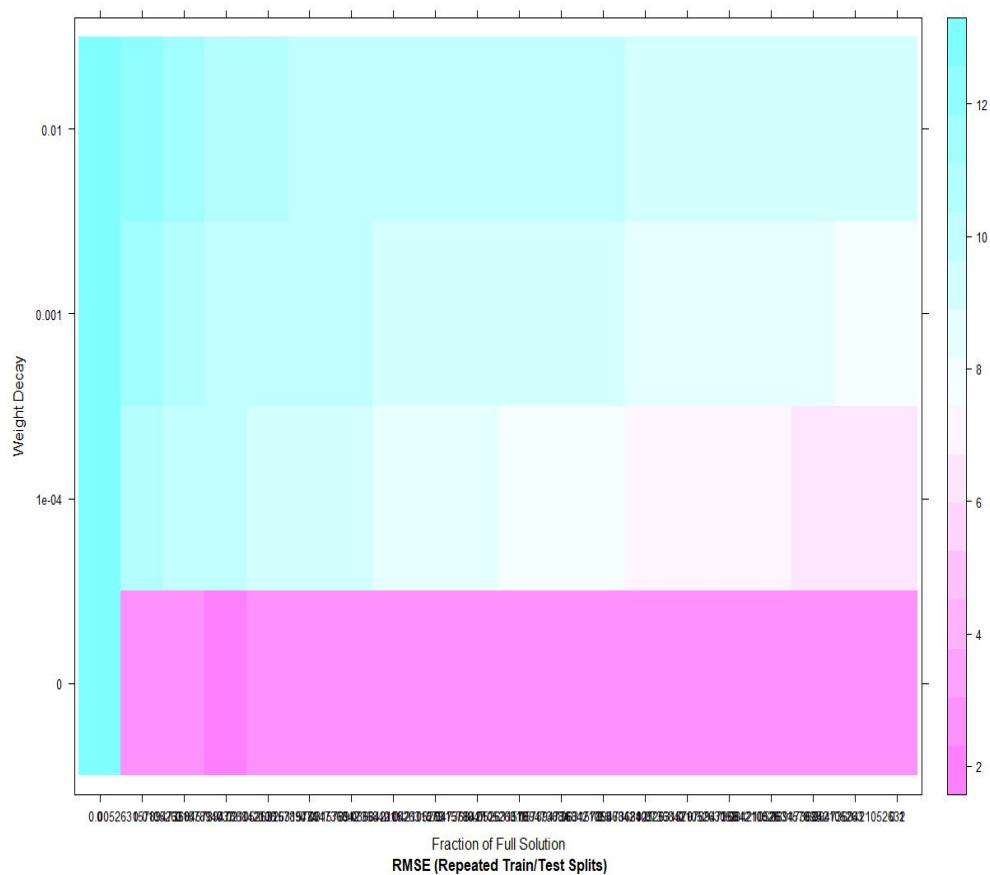
lambda	fraction	RMSE	Rsquared	MAE
oe+oo	0.000000000	12.574961	NaN	10.696237
oe+oo	0.005263158	2.526463	0.9622234	1.959691
oe+oo	0.010526316	2.335365	0.9682374	1.794543
oe+oo	0.015789474	2.294162	0.9698386	1.740583

RMSE was used to select the optimal model using the smallest value.

The final values used for the model were fraction = 0.01578947 and lambda = 0.

Since the Lambda value for this model is 0 we could consider this model to be a lasso model in practice.





The above representation shows the enet model performance for various parameters of lambda and fraction.

Best tuning parameters for each of the models are represented below:

	Ridge	Lasso	Enet
lambda	5.263158e-05	NA	0.00000000
fraction	NA	0.01578947	0.01578947

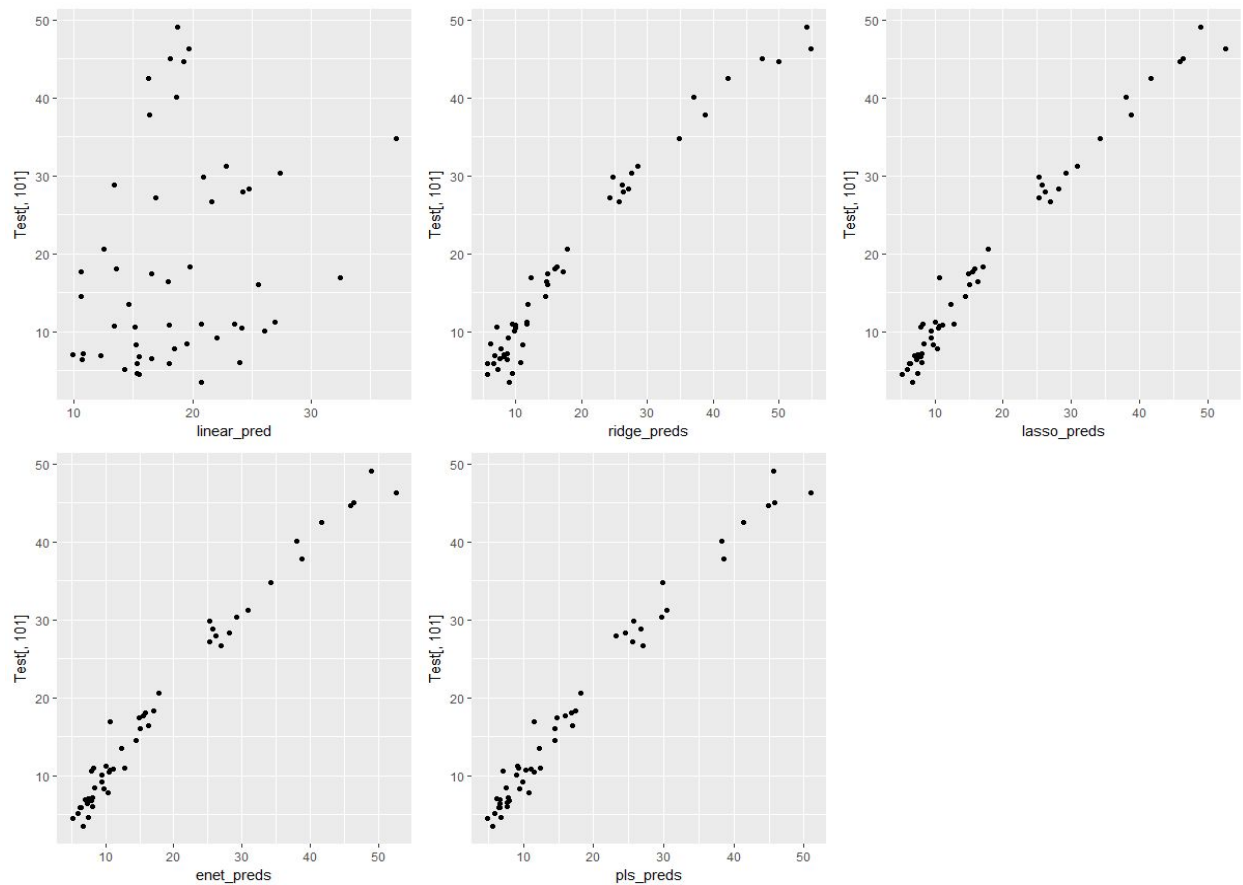
The table and graph for Enet agree with each other.

d) Which model has the best predictive ability? Is any model significantly better or worse than the others?

As the plot shows, the various models do a good job at predicting, except linear regression model.

Linear	PLS	Ridge	Lasso	Enet
0.2154802	0.9868296	0.9782466	0.9860839	0.9860839

The above table shows the values for the correlation between the predicted and observed values for each of the models. The Enet as well as the PLS model seems to be the best performing model amongst all the models. The PLS model seems to be the best model; this could be due to the nature of collinearity of the dataset. The collinearity may also be the reason for the worst performing model name i.e. linear regression.



(e) Explain which model you would use for predicting the fat content of a Sample

The plot shows the confidence intervals for the various models that we have implemented. To make the call between Enet and PLS we can refer to the R squared values for Elastic net it is 0.9698386 and for PLS it is 0.9649078.

