

Machine Learning Report

Red Wine Quality Dataset Analysis

by Mihai Calin

CEN4 S1.A

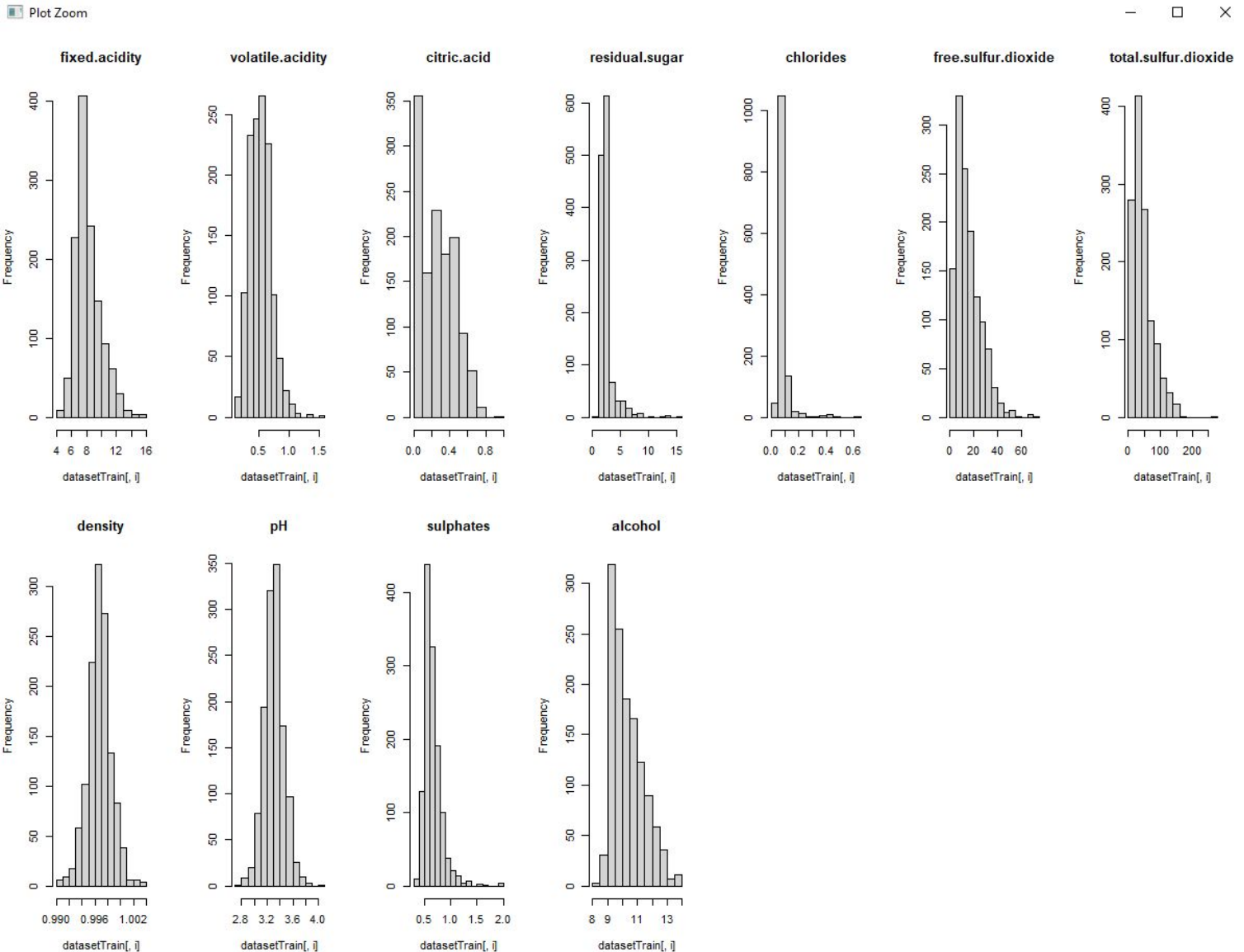
Red Wine Quality Dataset

Description and Initial Observations

The Red Wine Quality Dataset provides us with a selection of red wines and their certain characteristics, which values vary from wine to wine, and each wine has a certain quality.

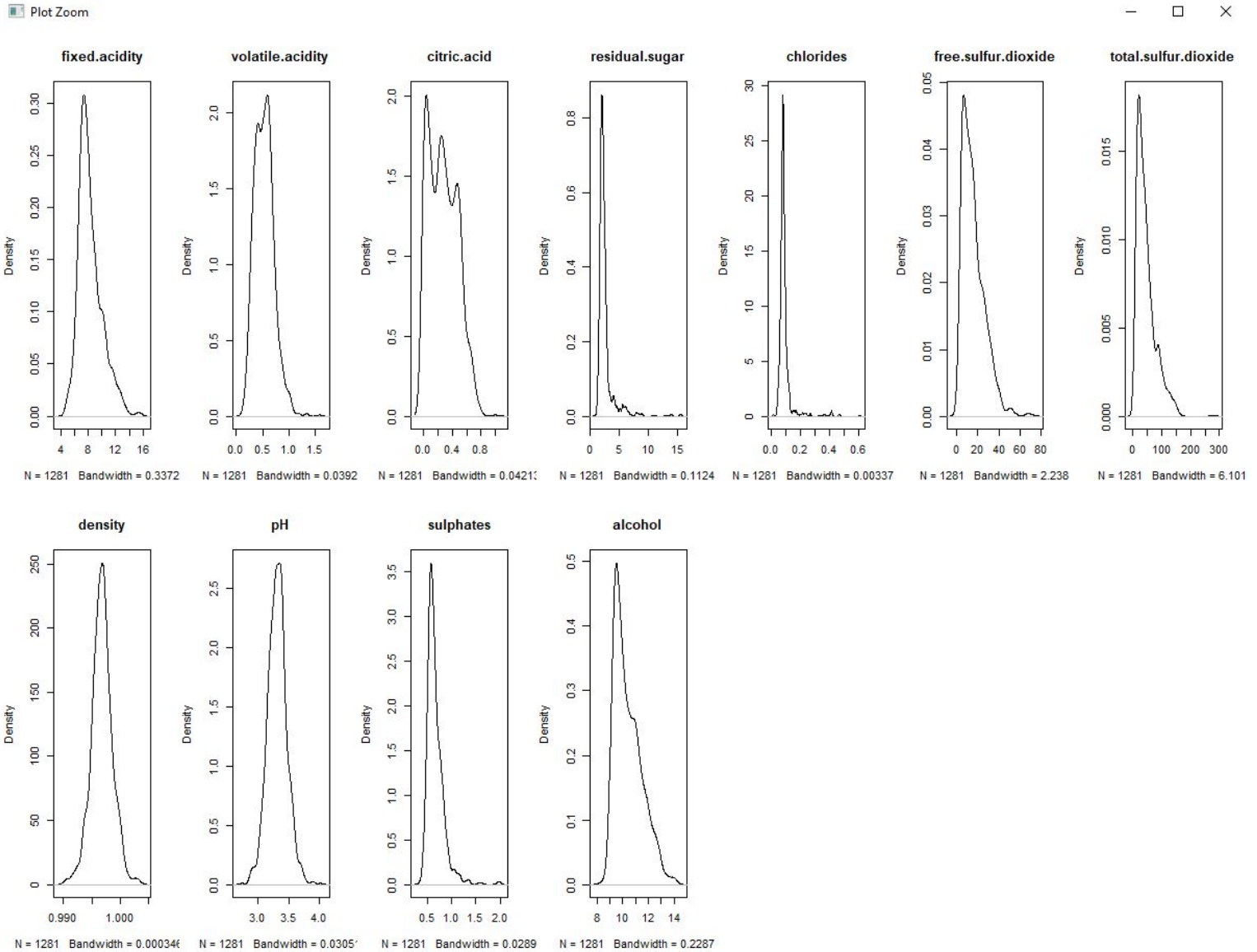
Based on the information available on the dataset, and using machine learning algorithms I'll try to train the data so that it'll predict the quality of the red wine with the least margin of error possible.

When it comes to the dataset, it's a reasonable large dataset, in which we have 12 variables (that make for the characteristics of our dear wines, ex: fixed acidity, volatile acidity, pH, residual sugar, quality etc).



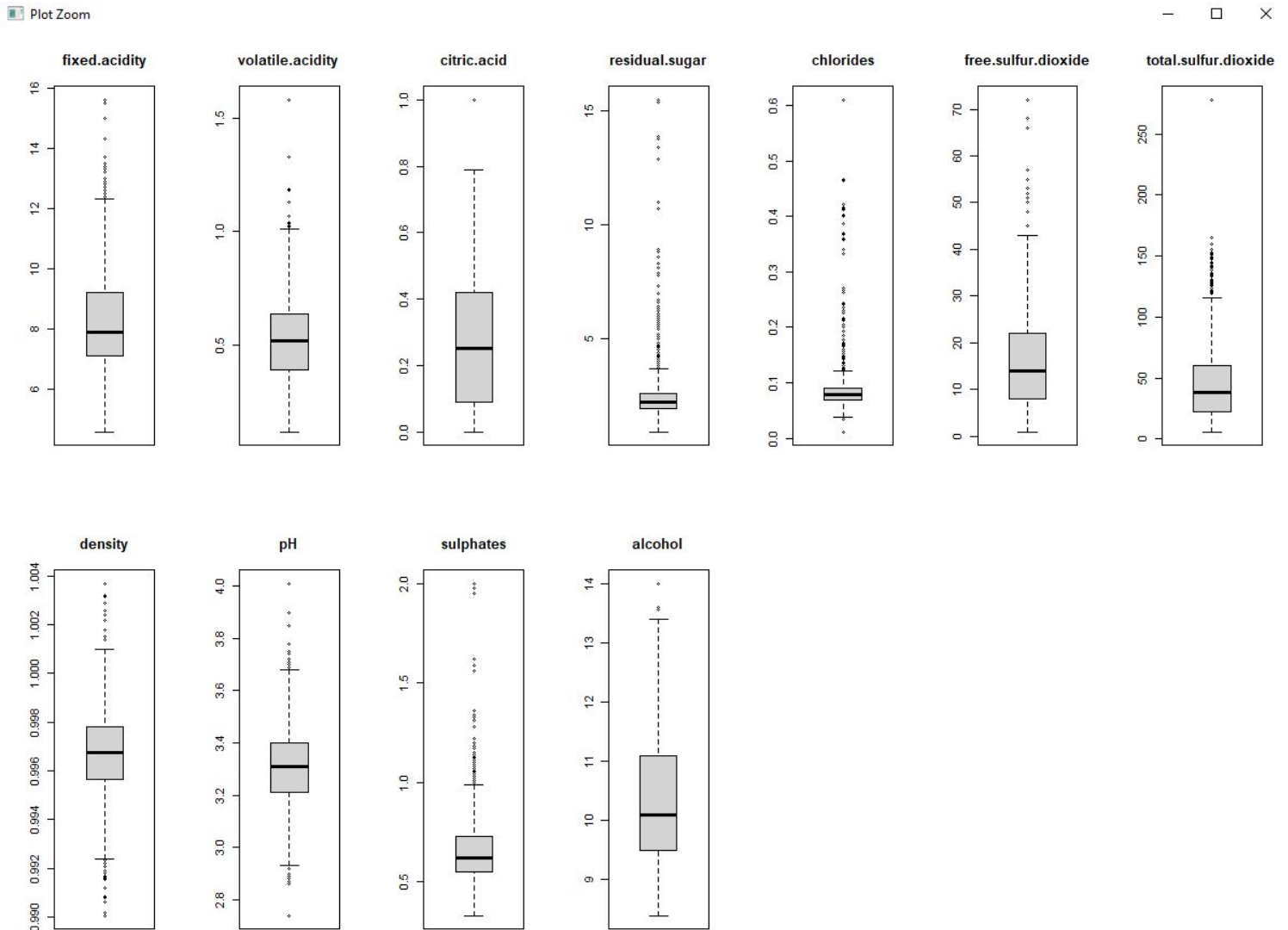
A short summary on the data:

- citric acid - seems to be somewhat uniformly distributed
- residual sugar - has a min - 0.9, and a max - 15, which is a massive difference
- chlorides - same as residual sugar, min - 0.012, max - 0.611
- free sulfur dioxide, total sulfur dioxide - same explanation as above



A short discussion on the results regarding the outliers :

- residual sugar: extreme outliers on the higher values
- chlorides: same as above
- sulphates: same as above, but not as extreme
- total sulphur dioxide: there is a concentration of outliers on the higher values

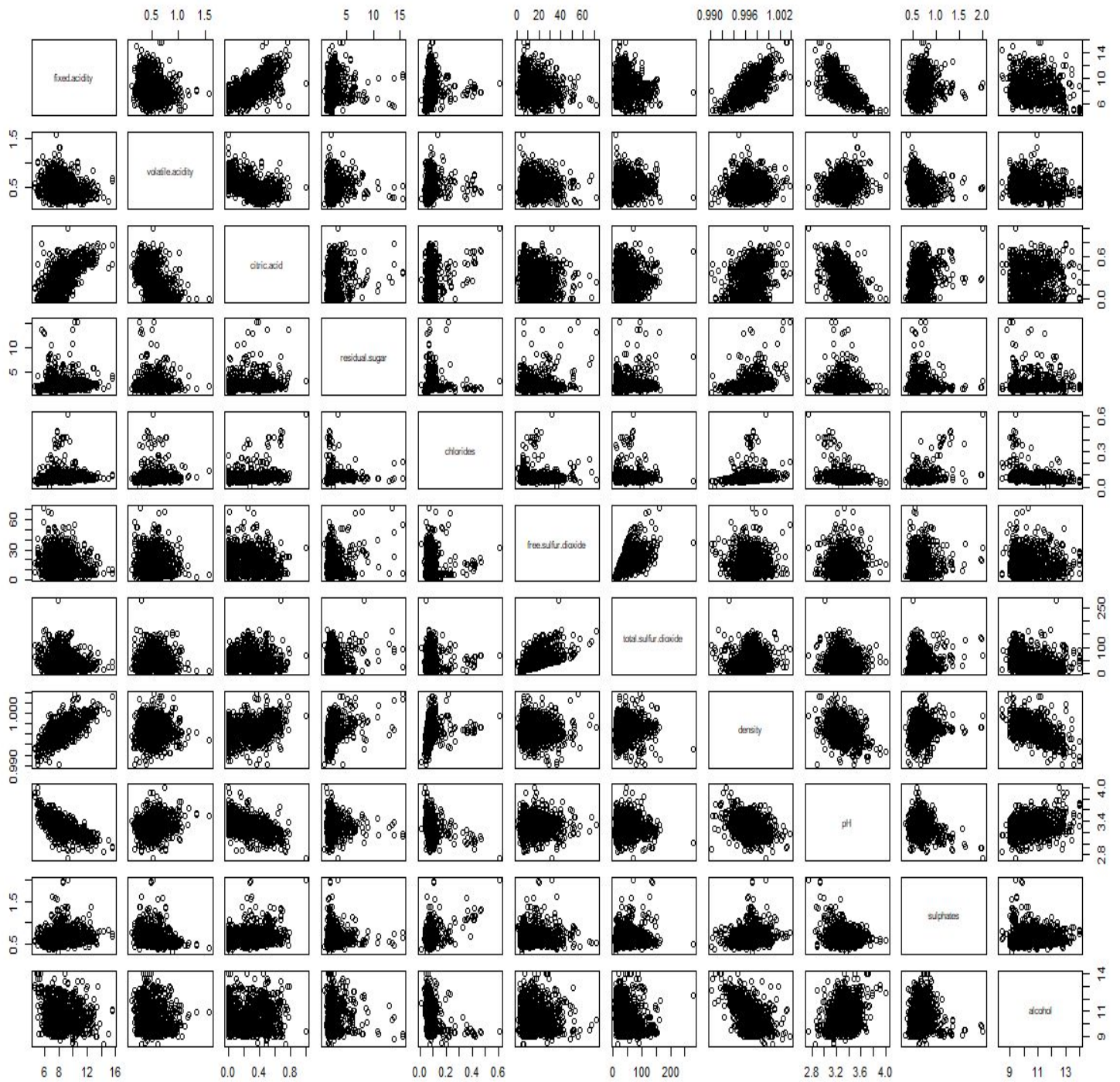


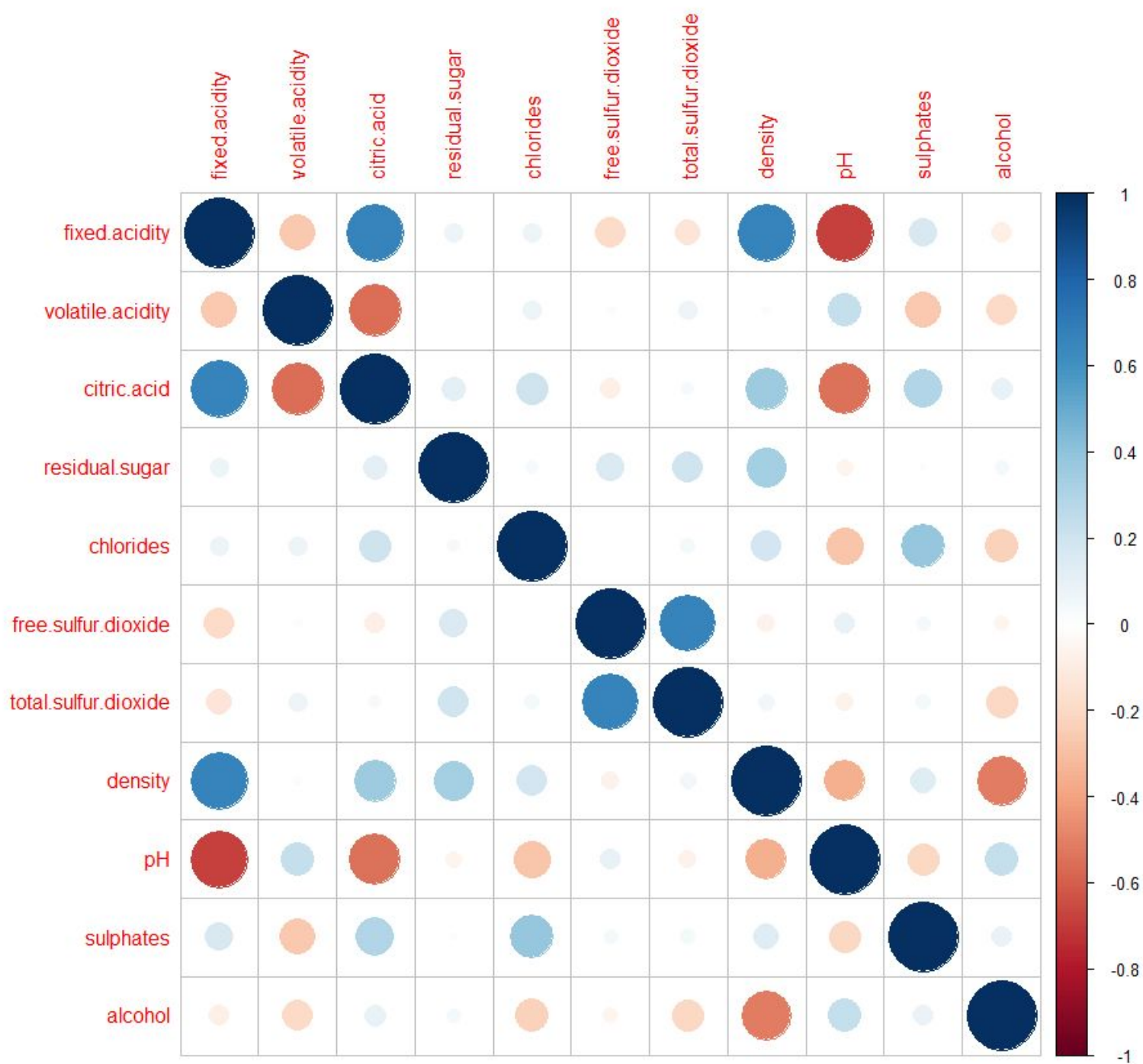
Since the excess of outliers is a known factor that impedes proper prediction I'll be trying to eliminate or at least reduce them, but for now we'll go ahead and see the initial results, so we can determine what impact they have on our predictions.

Now let's see the correlations:

Plot Zoom

— □ ×

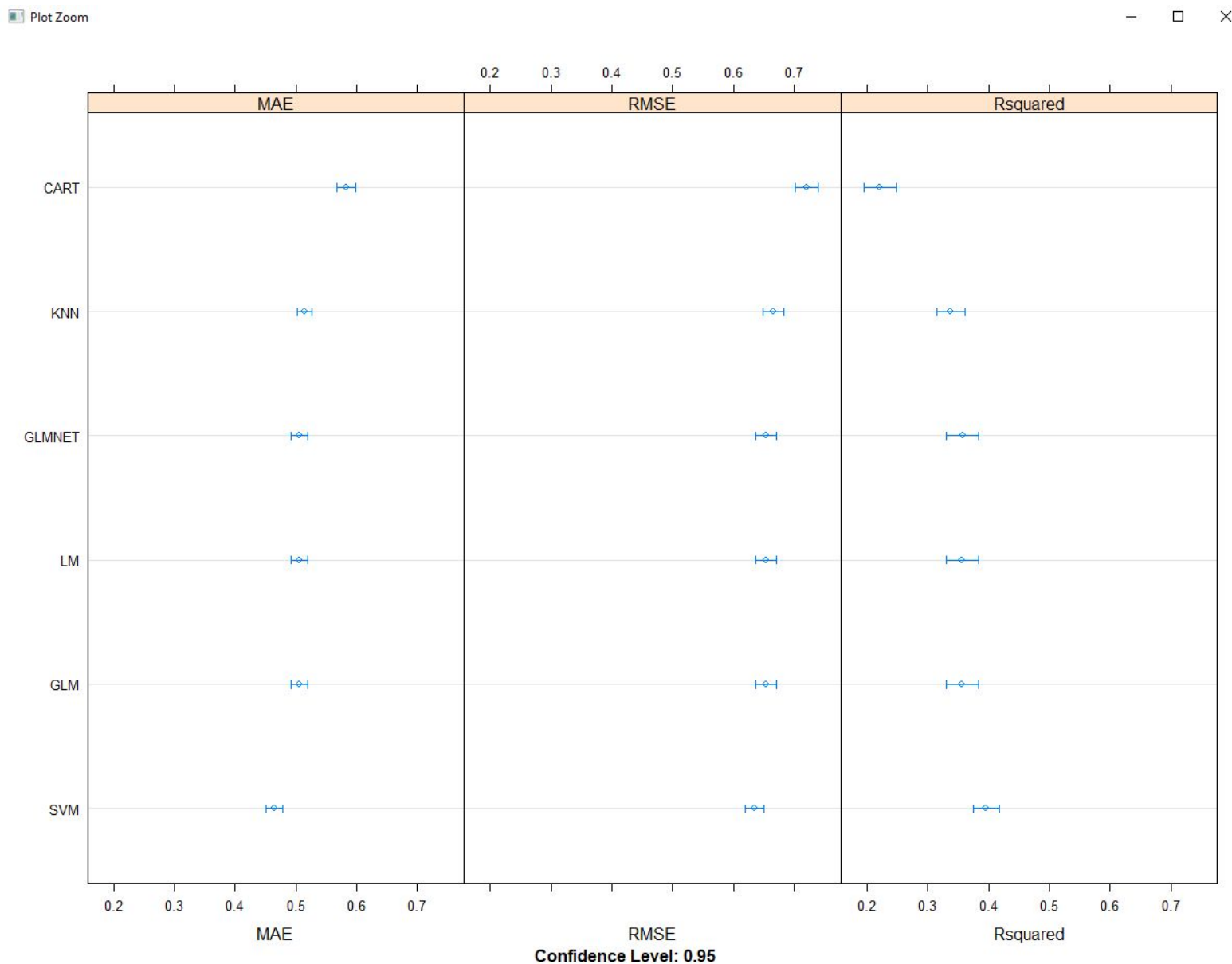




And the results are quite pleasant, because we don't have high correlations between variables, even though the fixed acidity and the free sulphur dioxide have a relatively higher correlation than the other variables.

Now to run the algorithms and begin the train of the data, I've split the dataset into 2 other datasets of different ration: one for training (80% of the data), and one for testing (20% of the data).

The result of the first run of the algorithms:



And the results seem promising, the margin of error is pretty small (RMSE), similarly with the median absolute error(MAE), which is good because smaller the two parameters are, then the precision of the predictions is going to be higher.

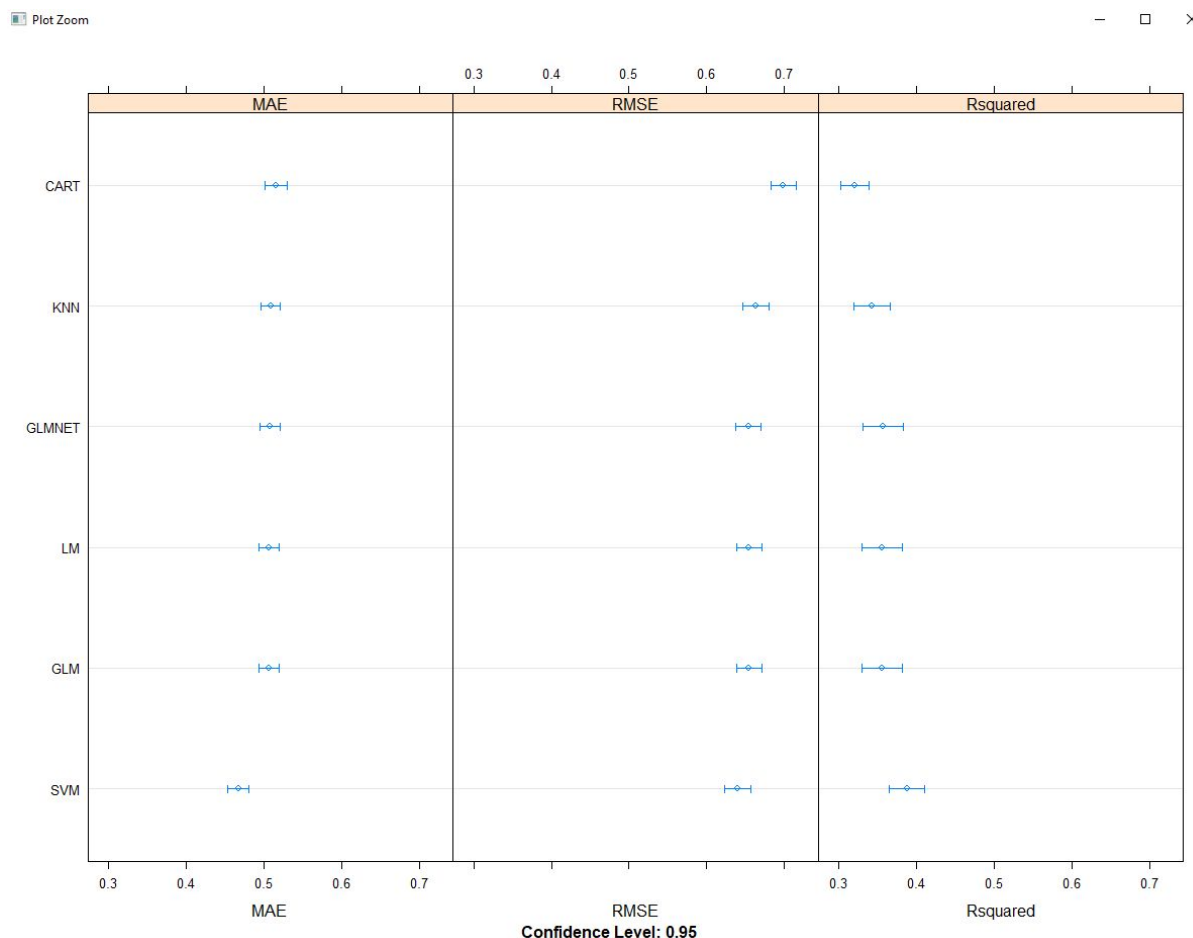
Rgarding the Rsquared, here the result is not so good since it's pretty close to 0, and the closer the value of Rsquared is to 1, the better the result.

Now I've tried to search for the variables that are highly correlated and eliminate them from the dataset, to see how that is going to affect the result.

I've slowly reduced the correlation factor for the cutoff, and only at 60% I found 2 variables that much correlation.

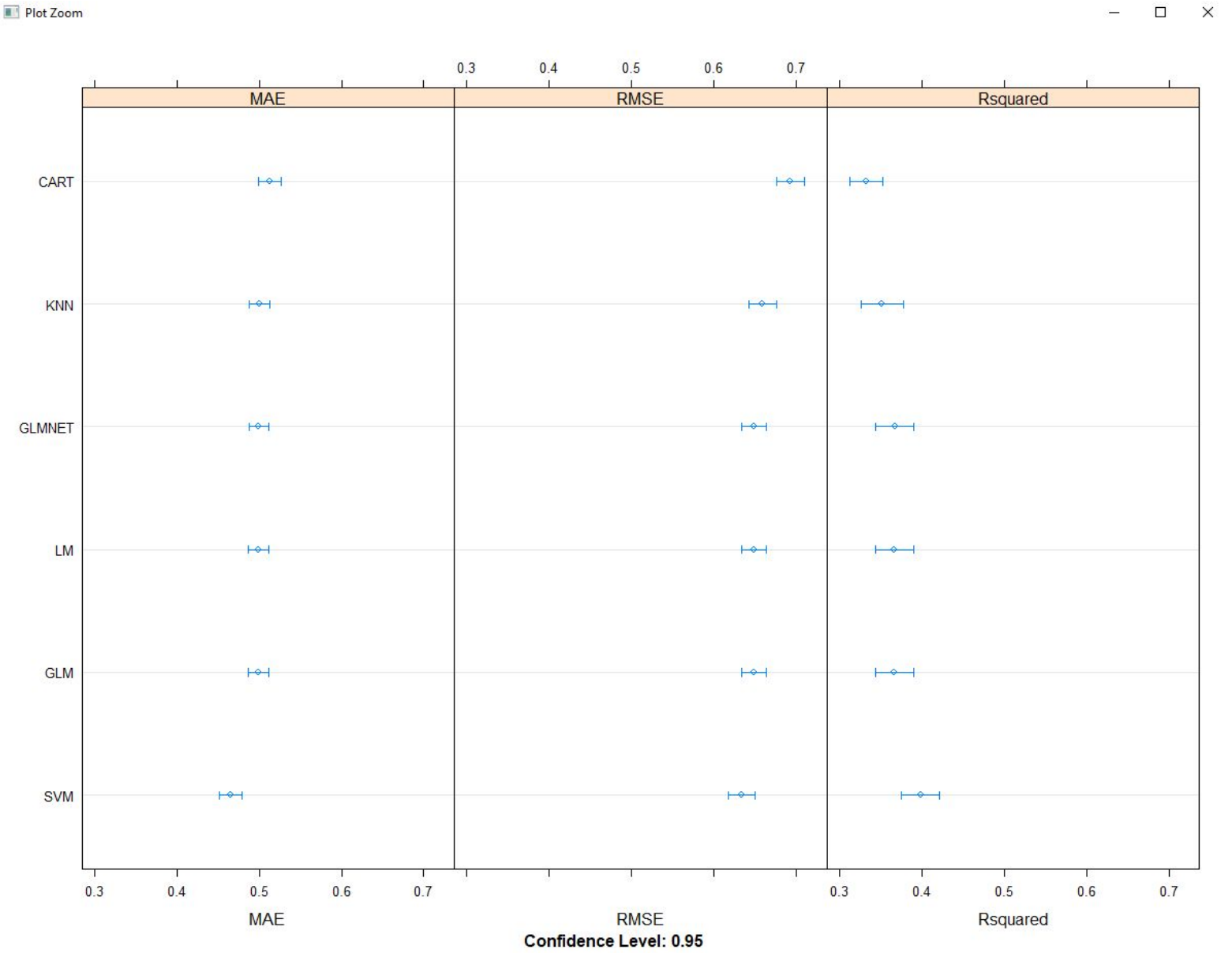
```
> # find attributes that are highly corrected
> set.seed(7)
> cutoff <- 0.60
> correlations <- cor(datasetTrain[,1:10])
> highlyCorrelated <- findCorrelation(correlations, cutoff=cutoff)
> for (value in highlyCorrelated) {
+   print(names(datasetTrain)[value])
+ }
[1] "fixed.acidity"
[1] "free.sulfur.dioxide"
> # create a new dataset without highly corrected features
> dataset_features <- datasetTrain[, -highlyCorrelated]
> dim(dataset_features)
[1] 1361    9
>
```

Now, the results of the algorithms on the new dataset that doesn't include the highly correlated variables:



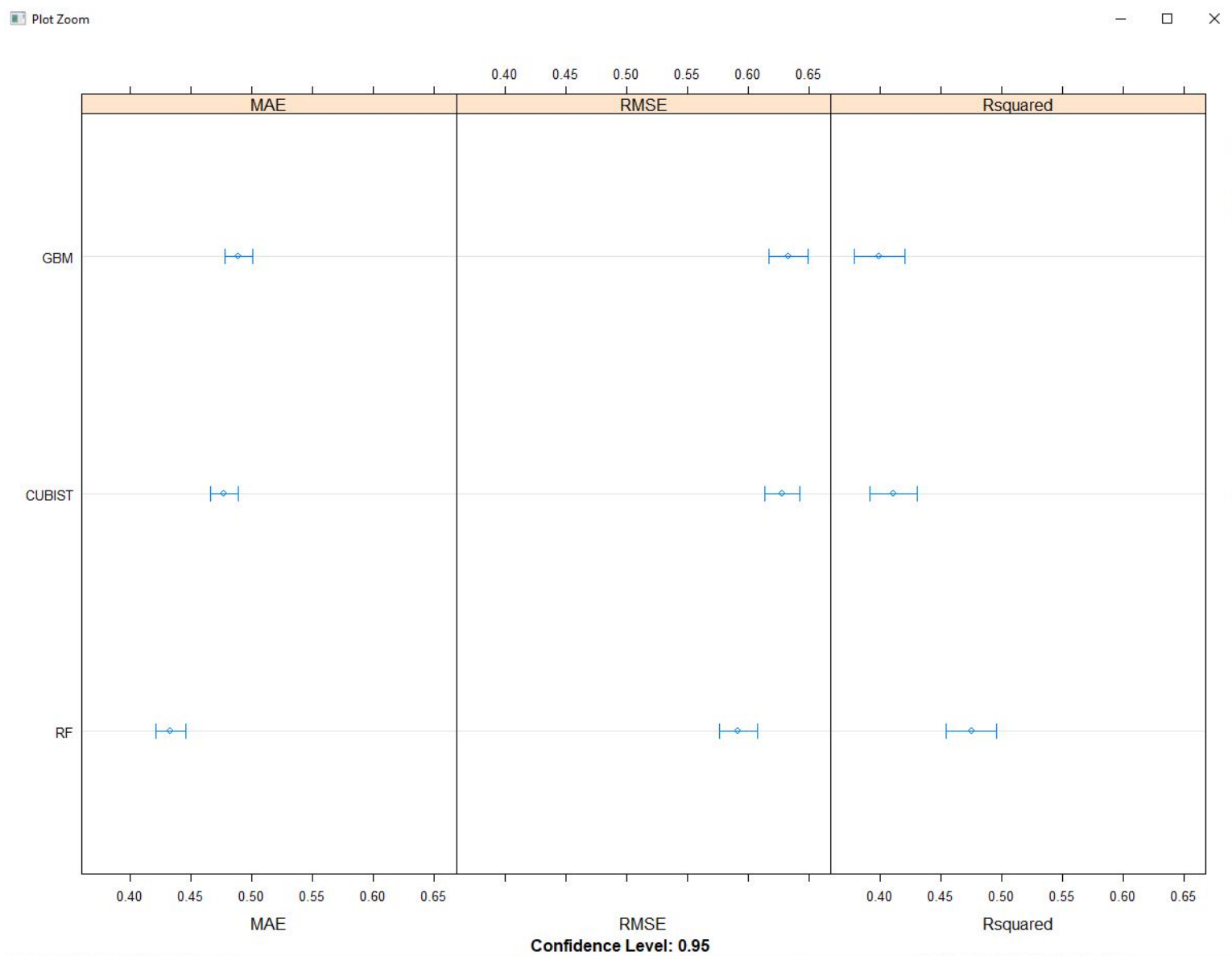
And the changes are insignificant, so we'll proceed with the initial dataset.

Now, let's try to run the algorithms with the Box-Cox Transform:



And, again the changes proved insignificant.

Now, let's try to run other algorithms, to see if it'll have any significant changes:



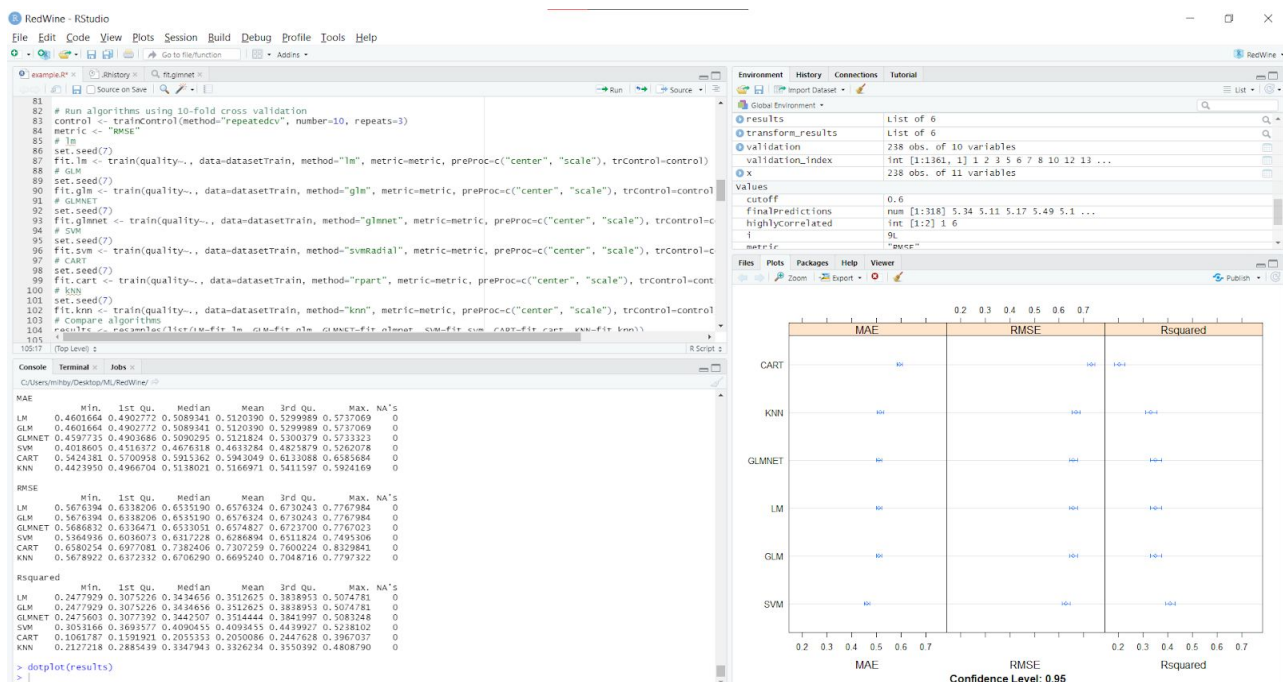
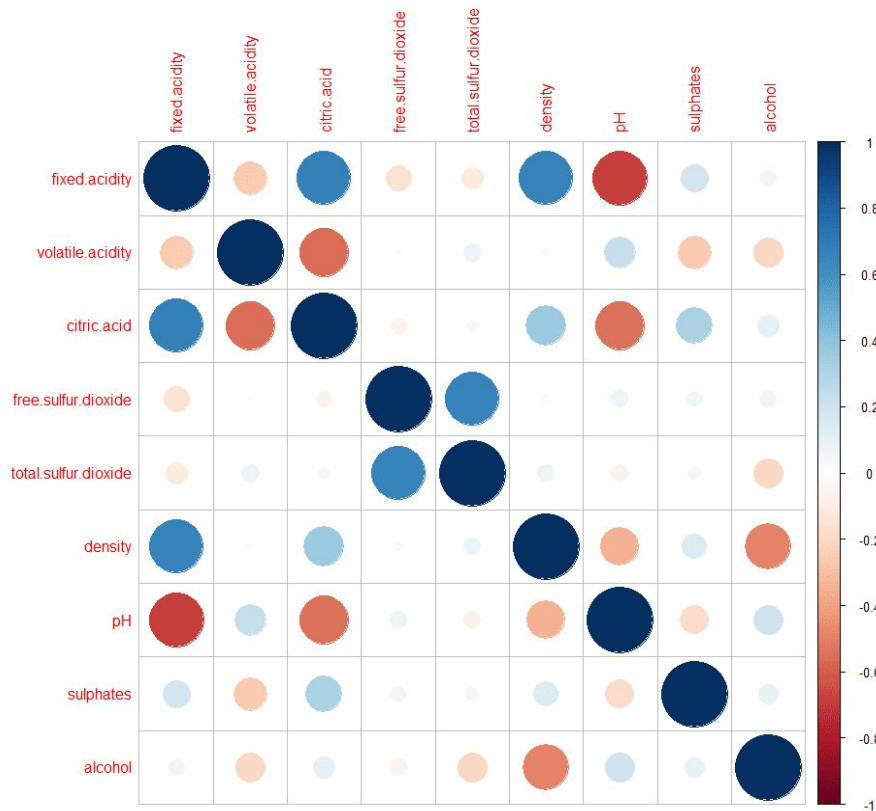
And now we have a better result, so I'll be using the this algorithm to run on new data (the remaining dataset that I've saved for testing) and see how far off are the predictions:

```

Console Terminal Jobs
C:/Users/mibby/Desktop/ML/RedWine/
>
> predictions <- predict(fit.cubist, newdata=x)
> print(predictions)
[1] 5.518986 5.691974 5.258907 4.948594 5.516471 5.342236 5.240856 4.802408 5.385659 5.384938 4.571118 5.127559 5.172316
[14] 5.165021 5.395542 5.224705 5.122672 5.435228 5.607400 5.896239 4.995122 5.265761 5.704646 5.349718 4.941485 5.225691
[27] 5.045173 4.942307 5.601873 4.527615 6.212644 5.421463 5.055743 5.314455 5.099916 5.445649 5.182467 5.366728 6.214897
[40] 5.546139 5.433943 5.107702 5.372950 5.764687 6.137141 5.743628 6.372781 6.016290 6.448633 5.139544 6.269341 6.542992
[53] 5.933819 5.337881 5.140034 5.250845 5.886510 6.060398 5.378976 6.232417 5.611037 6.386326 5.808299 5.756944 6.854245
[66] 6.364571 5.185845 6.271636 5.158449 5.744769 6.053005 5.618264 5.694557 5.462574 5.739222 5.626937 5.657332 5.451989
[79] 4.944979 6.763140 5.429212 5.384243 5.388945 5.440673 4.938229 5.144464 5.008227 6.855949 4.679731 5.357247 5.367270
[92] 4.711611 5.711736 6.004842 5.364109 4.575601 4.586238 5.535389 5.742434 5.134079 4.963684 5.109643 5.053529 5.151733
[105] 5.042246 5.544608 6.020737 5.200624 5.666564 5.293425 7.066842 5.601353 5.724722 5.718869 5.281573 5.565941 6.446341
[118] 4.906415 4.976942 5.507220 5.004480 5.268829 5.253993 6.193238 6.965869 6.633873 5.049393 5.049393 5.246224 5.741293
[131] 6.351403 6.467686 5.623058 6.341032 6.901127 5.335862 6.492669 5.518254 6.462317 6.192672 5.494213 4.775688 5.591523
[144] 4.976056 6.883724 5.203425 6.357997 5.325560 5.084665 5.601664 5.619888 6.190652 7.103249 5.975210 6.334642 6.761431
[157] 5.327249 6.201297 5.991802 6.751660 5.903841 5.903841 6.060546 6.319194 7.059261 6.052242 6.052242 5.778294 5.869025
[170] 5.970800 6.430749 5.747671 5.650014 5.398863 5.409987 5.269711 5.444095 6.487516 6.489898 6.758188 6.081740 5.500779
[183] 6.480591 5.992816 5.340658 5.842066 6.936607 5.543820 4.967350 6.107648 6.172801 4.975404 5.259758 5.591878 5.636509
[196] 5.636509 6.050605 5.694395 4.736331 5.150850 5.252522 5.907645 5.848881 4.966335 5.123662 5.126166 5.372354 5.356155
[209] 4.998581 5.112996 5.422979 6.715468 5.757374 5.298185 5.554851 6.111851 5.748778 5.583792 5.532259 5.467219 5.658670
[222] 5.441863 4.787041 5.267632 6.070026 5.959546 5.322480 5.735471 5.577078 4.993229 6.263361 6.144412 6.252946 5.679534
[235] 7.064692 6.354968 6.361186 6.032565
>
> # calculate RMSE
> rmse <- RMSE(predictions, y)
> r2 <- R2(predictions, y)
> print(rmse)
[1] 0.615174
>

```

To improve the results I've tried to eliminate the outliers, towards this goal I've removed the variables with the most outliers (chlorides and residual sugar) and compare the results:



And again, the changes are insignificant, we'll proceed with the initial dataset.

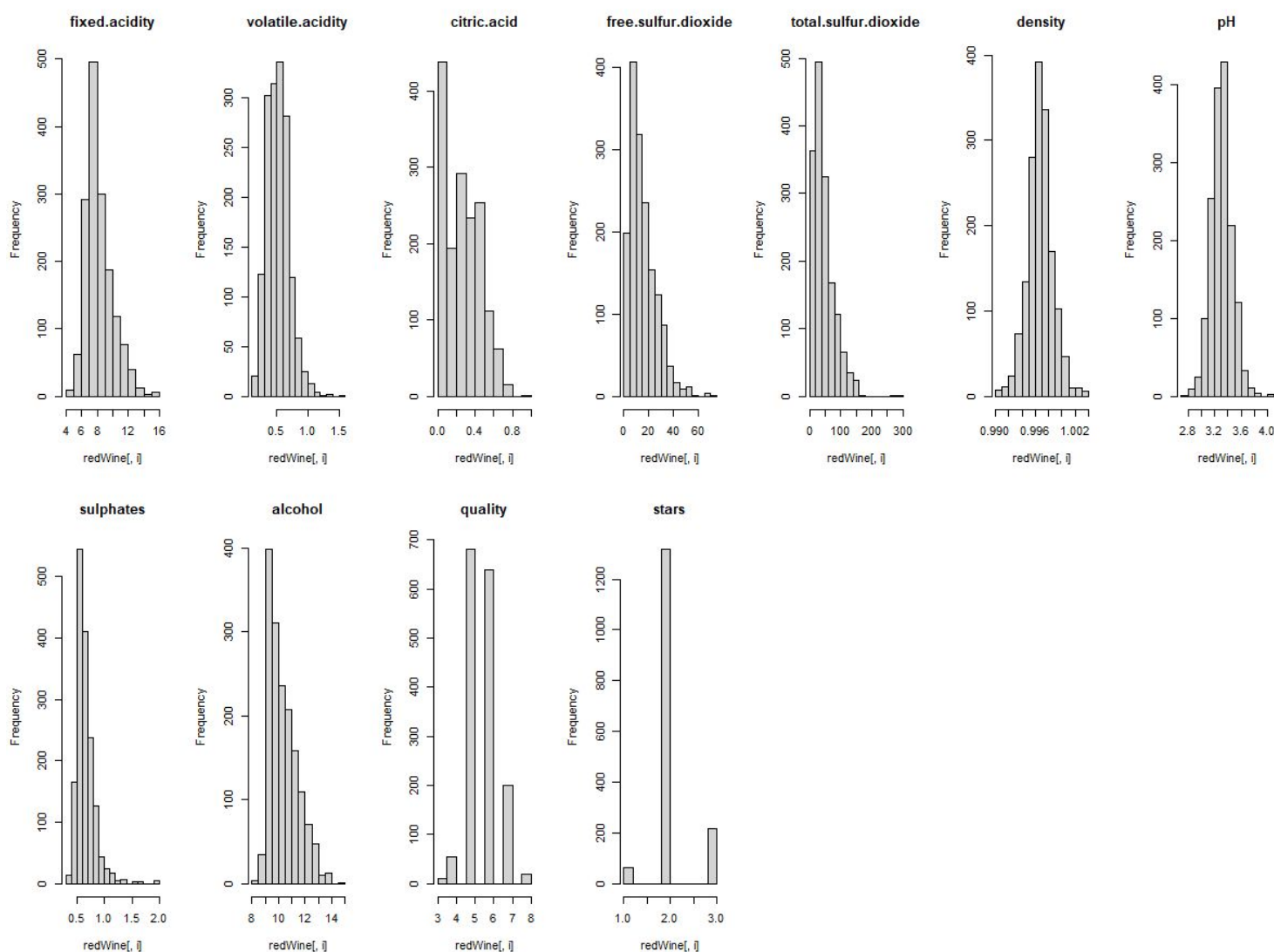
Hoping that I'll improve the results I've changed the split ration of the two datasets (85/15) leaving more data for the training dataset.

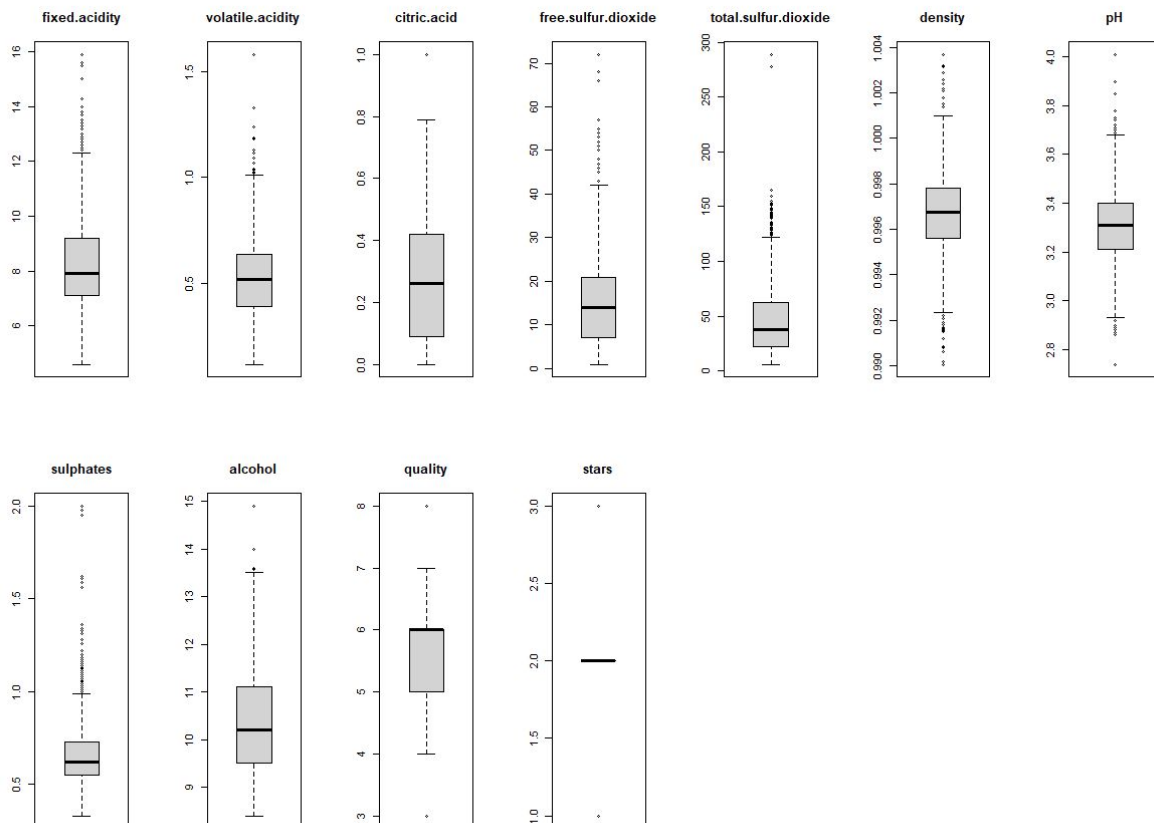
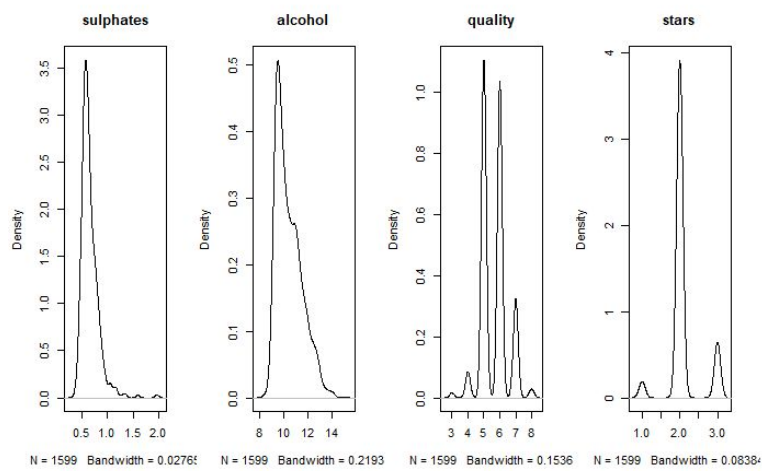
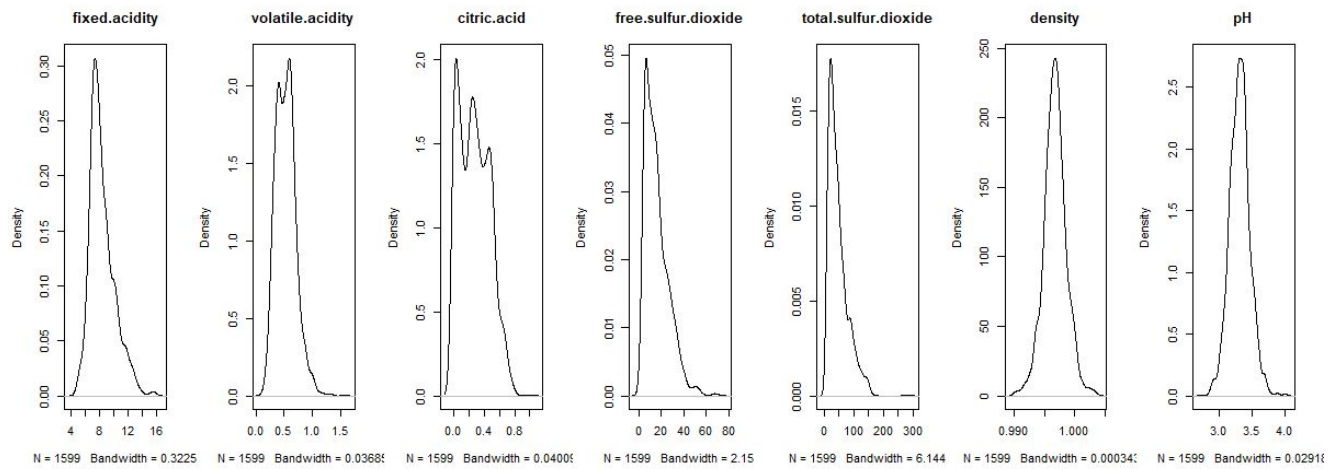
And I've changed the output variable, creating a new system of rating the quality of the red wines, the new variable is names "Stars" and has 3 possible values:

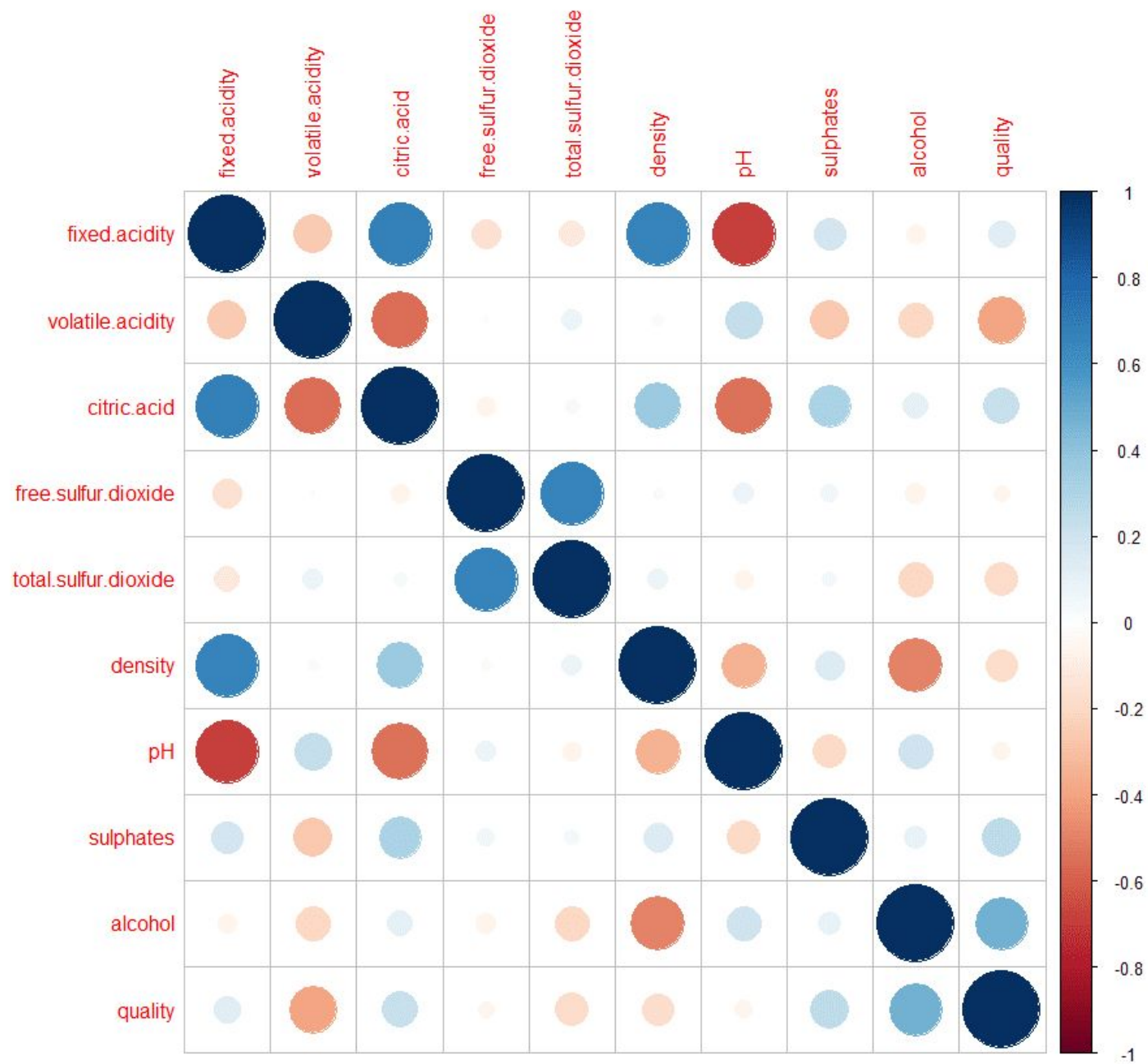
- 1 is for the wines that have a quality below 5 (poor quality)
- 2 is for the wines that have a quality between 5 and 7 (good quality)
- 3 is for the wines that have a quality over 7 (great quality)

And with this addition we now have the "Quality" as an input variable that should provide valuable information.

And now I'll repeat the process and evaluate the changes.








```

101 corplot(correlations, method="circle")
102
103 # Evaluate Algorithms
104
105 # Run algorithms using 10-fold cross validation
106 control <- trainControl(method="repeatedcv", number=10, repeats=3)
107 metric <- "RMSE"
108 # lm
109 set.seed(7)
110 fit.lm <- train(stars~, data=datasetTrain, method="lm", metric=metric, preProc=c("center", "scale"), trControl=control)
111 # GLM
112 set.seed(7)
113 fit.glm <- train(stars~, data=datasetTrain, method="glm", metric=metric, preProc=c("center", "scale"), trControl=control)
114 # GLMNET
115 set.seed(7)
116 fit.glmnet <- train(stars~, data=datasetTrain, method="glmnet", metric=metric, preProc=c("center", "scale"), trControl=control)
117 # SVM
118 set.seed(7)
119 fit.svm <- train(stars~, data=datasetTrain, method="svmRadial", metric=metric, preProc=c("center", "scale"), trControl=control)
120 # CART
121 set.seed(7)
122 fit.cart <- train(stars~, data=datasetTrain, method="rpart", metric=metric, preProc=c("center", "scale"), trControl=control)
123 # KNN
124 set.seed(7)
125 fit.knn <- train(stars~, data=datasetTrain, method="knn", metric=metric, preProc=c("center", "scale"), trControl=control)
126 # Compare algorithms
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

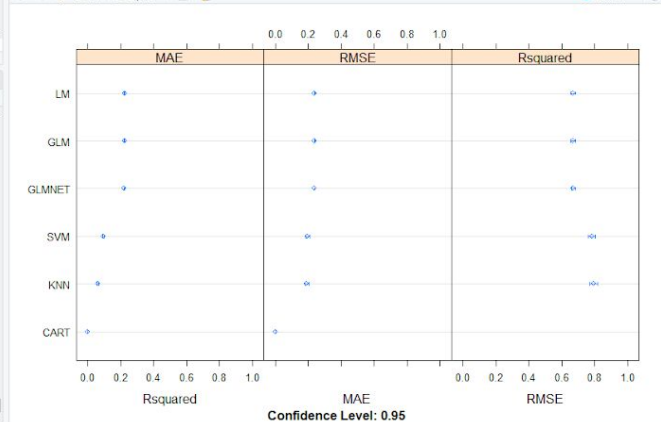
```

MAE	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
LM	0.21291910	0.22196258	0.22413786	0.22435863	0.22761482	0.23429523	0
GLM	0.21291910	0.22196258	0.22413786	0.22435863	0.22761482	0.23429523	0
GLMNET	0.21281809	0.22114892	0.22345240	0.22334189	0.22631103	0.23162422	0
SVM	0.06091681	0.08861527	0.09793761	0.09563448	0.10377077	0.11466709	0
CART	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0
KNN	0.02481752	0.05474887	0.06446078	0.06197422	0.06910755	0.09632353	0

RMSE	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
LM	0.22754228	0.2347780	0.2373660	0.2380913	0.2418473	0.2504369	0
GLM	0.22754228	0.2347780	0.2373660	0.2380913	0.2418473	0.2504369	0
GLMNET	0.22796763	0.2335905	0.2375732	0.2372187	0.2408964	0.2473214	0
SVM	0.10974502	0.1884968	0.2032667	0.1959726	0.2164457	0.2434713	0
CART	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0
KNN	0.09201714	0.1773407	0.1937769	0.1901744	0.2078232	0.2767842	0

Rsqured	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
LM	0.5733101	0.6457981	0.6652919	0.6673890	0.6946214	0.7235360	0
GLM	0.5733101	0.6457981	0.6652919	0.6673890	0.6946214	0.7235360	0
GLMNET	0.5805004	0.6475710	0.6693372	0.6701422	0.6942450	0.7227298	0
SVM	0.6631808	0.7531028	0.7754566	0.7831887	0.8002890	0.9241230	0
CART	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	0
KNN	0.5789519	0.7611453	0.8029461	0.7938939	0.8238951	0.9500757	0

Environment	History	Connections	Tutorial
Global Environment			
surpinales	num 0.36 0.68 0.65 0.38 0.36 0.36 0.46 0.47 0.37 0.8 ...		
alcohol	num 9.4 9.8 9.8 9.8 9.4 9.4 10 9.5 10.5 ...		
quality	int 5 5 5 6 5 5 7 7 5 ...		
stars	num 2 2 2 2 2 2 2 3 2 ...		
results	List of 6		
transform_results	List of 6		
validation	238 obs. of 11 variables		
validation_index	int [1:1361, 1] 1 2 3 5 6 7 8 10 12 13 ...		
x	238 obs. of 9 variables		
values			
cutoff	0.6		
finalPredictions	num [1:238] 5.52 5.69 5.26 5.26 4.95 5.52 ...		
highlyCorrelated	int [1:2] 1 4		
i	11L		
metric	"RMSE"		



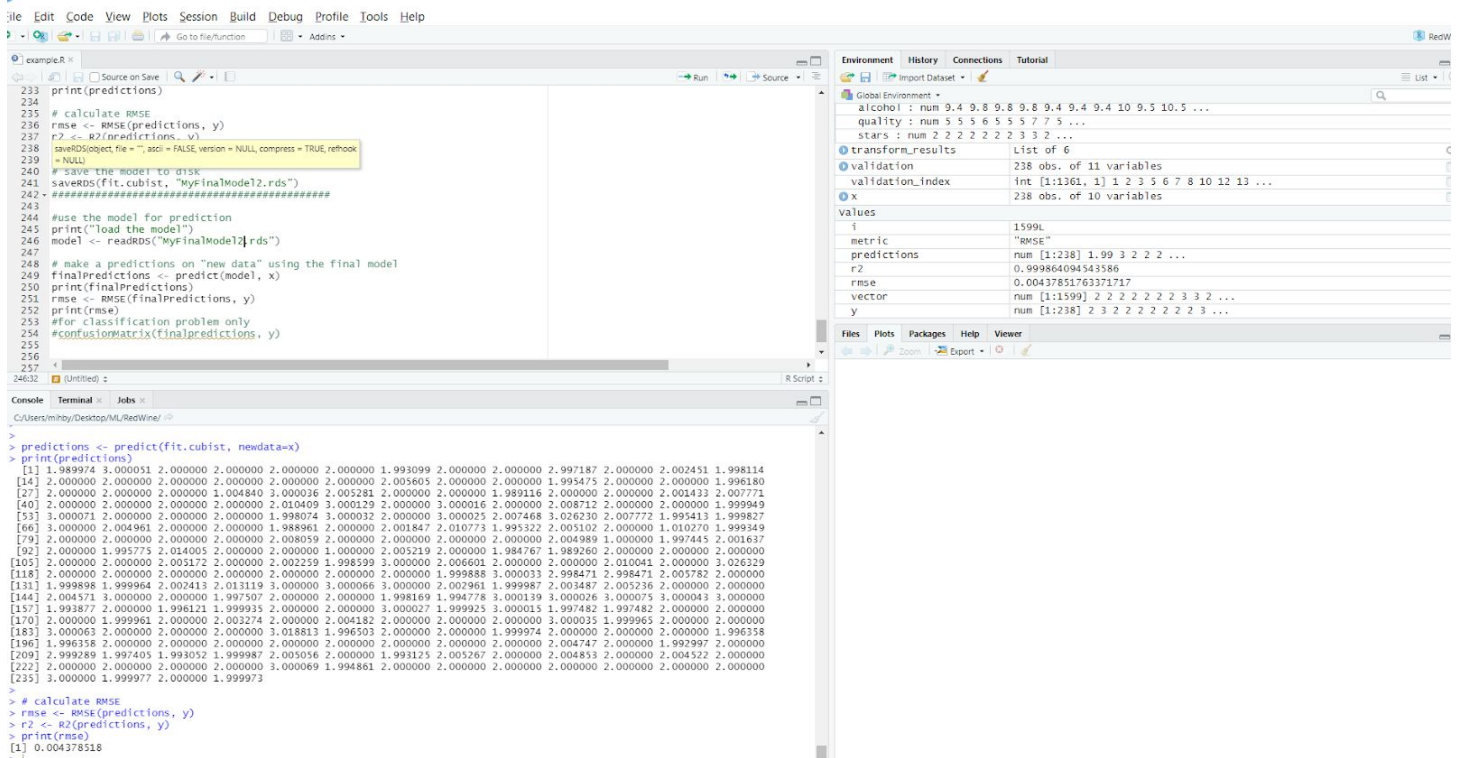
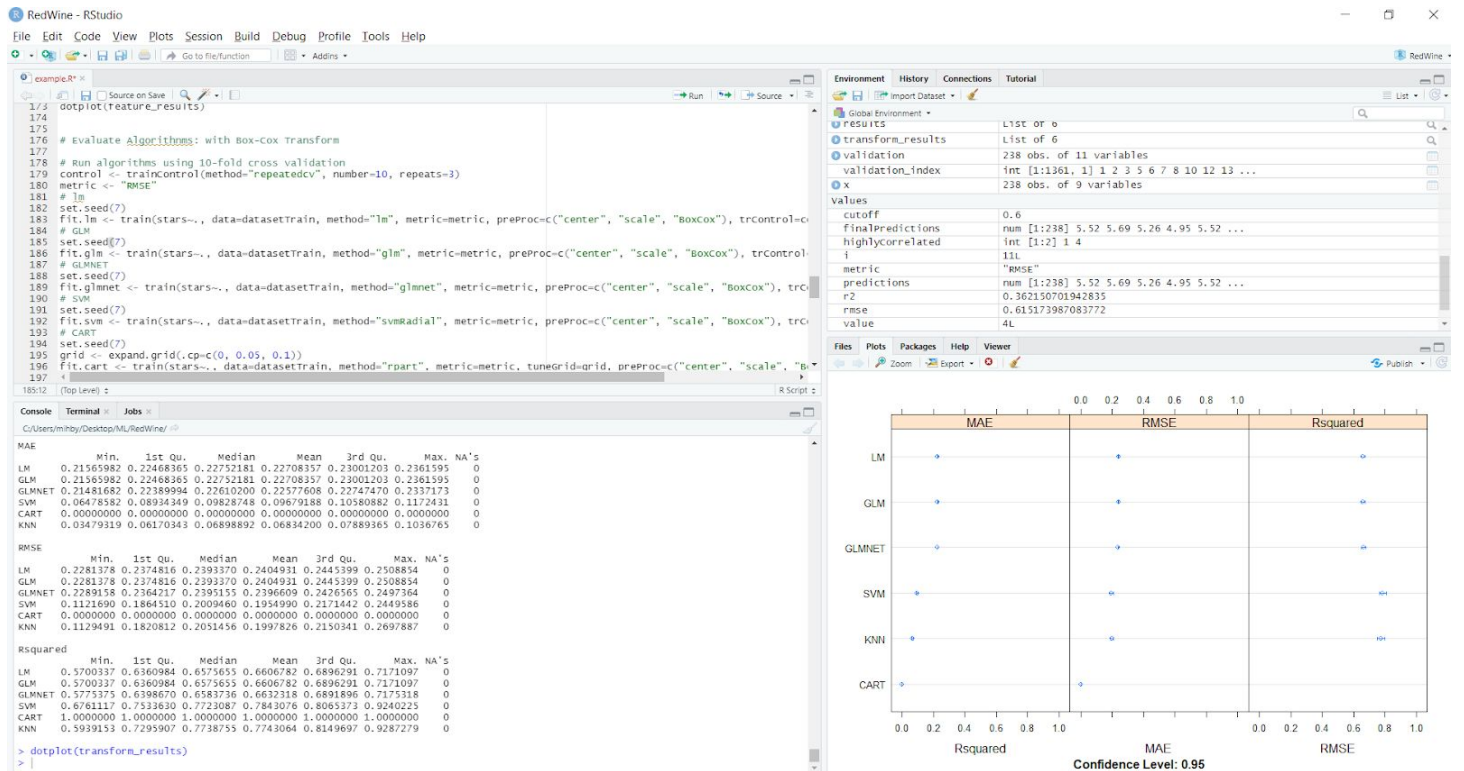
```

134 # Evaluate Algorithms: With Feature Selection Step
135
136 # remove correlated attributes
137 # find attributes that are highly correlated
138 set.seed(7)
139 cutoff <- 0.60
140 correlations <- cor(datasetTrain[,1:10])
141 highlyCorrelated <- findCorrelation(correlations, cutoff=cutoff)
142 for (value in highlyCorrelated) {
143   print(names(datasetTrain)[value])
144 }
145 # create a new dataset without highly correlated features
146 dataset_features <- datasetTrain[,!highlyCorrelated]
147 dim(dataset_features)
148
149 # Run algorithms using 10-fold cross validation
150 control <- trainControl(method="repeatedcv", number=10, repeats=3)
151 metric <- "RMSE"
152 # lm
153 set.seed(7)
154 fit.lm <- train(stars~, data=dataset_features, method="lm", metric=metric, preProc=c("center", "scale"), trControl=control)
155 # GLM
156 set.seed(7)
157 fit.glm <- train(stars~, data=dataset_features, method="glm", metric=metric, preProc=c("center", "scale"), trControl=control)
158 # GLMNET
159 set.seed(7)
160 fit.glmnet <- train(stars~, data=dataset_features, method="glmnet", metric=metric, preProc=c("center", "scale"), trControl=control)
161 # SVM
162 set.seed(7)
163 fit.svm <- train(stars~, data=dataset_features, method="svmRadial", metric=metric, preProc=c("center", "scale"), trControl=control)
164 # CART
165 set.seed(7)
166 fit.cart <- train(stars~, data=dataset_features, method="rpart", metric=metric, preProc=c("center", "scale"), trControl=control)
167 # KNN
168 set.seed(7)
169 fit.knn <- train(stars~, data=dataset_features, method="knn", metric=metric, preProc=c("center", "scale"), trControl=control)
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

```

134 # Evaluate Algorithms: With Feature Selection Step
135
136 # remove correlated attributes
137 # find attributes that are highly correlated
138 set.seed(7)
139 cutoff <- 0.60
140 correlations <- cor(datasetTrain[,1:10])
141 highlyCorrelated <- findCorrelation(correlations, cutoff=cutoff)
142 for (value in highlyCorrelated) {
143   print(names(datasetTrain)[value])
144 }
145 # create a new dataset without highly correlated features
146 dataset_features <- datasetTrain[,!highlyCorrelated]
147 dim(dataset_features)
148
149 # Run algorithms using 10-fold cross validation
150 control <- trainControl(method="repeatedcv", number=10, repeats=3)
151 metric <- "RMSE"
152 # lm
153 set.seed(7)
154 fit.lm <- train(stars~, data=dataset_features, method="lm", metric=metric, preProc=c("center", "scale"), trControl=control)
155 # GLM
156 set.seed(7)
157 fit.glm <- train(stars~, data=dataset_features, method="glm", metric=metric, preProc=c("center", "scale"), trControl=control)
158 # GLMNET
159 set.seed(7)
160 fit.glmnet <- train(stars~, data=dataset_features, method="glmnet", metric=metric, preProc=c("center", "scale"), trControl=control)
161 # SVM
162 set.seed(7)
163 fit.svm <- train(stars~, data=dataset_features, method="svmRadial", metric=metric, preProc=c("center", "scale"), trControl=control)
164 # CART
165 set.seed(7)
166 fit.cart <- train(stars~, data=dataset_features, method="rpart", metric=metric, preProc=c("center", "scale"), trControl=control)
167 # KNN
168 set.seed(7)
169 fit.knn <- train(stars~, data=dataset_features, method="knn", metric=metric, preProc=c("center", "scale"), trControl=control)
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
```

So, when I'm running the algorithms to train the data, the mean RMSE is 0.2 which is a massive improvement from 0.6, and when I'm doing the predictions on the testing dataset the RMSE is 0.004 which is so low that makes our predictions too close to being perfect, which could mean that the dataset is overtrained, and I'll need new data to test if that is the actual RMSE.

But, overall massive improvement to the results due to the addition of "Stars" as an output variable.