**Assignment #6**
Nate Bitting

## Introduction

The purpose of this assignment is to explore the use of principal component analysis (PCA) to alleviate multicollinearity that may exist among the predictor variables when building a multiple regression model.  The dataset we will be working with is daily stock prices for twenty stocks (predictor variable candidates) included in a fund managed by Vanguard (VV).  The first step in our analysis will be to perform a log transformation on the daily returns (current – previous stock prices) for all of the individual stocks in the fund as well as the market index (VV), which we will use as our response variable. The second step in our analysis, after performing the transformations on the predictor and response variables, is to plot a Pearson correlation matrix for all of the predictor variables against the response variable, which we'll call response_VV. This analysis will identify any correlations between the predictor and response variables and identify any multicollinearity among the variables.
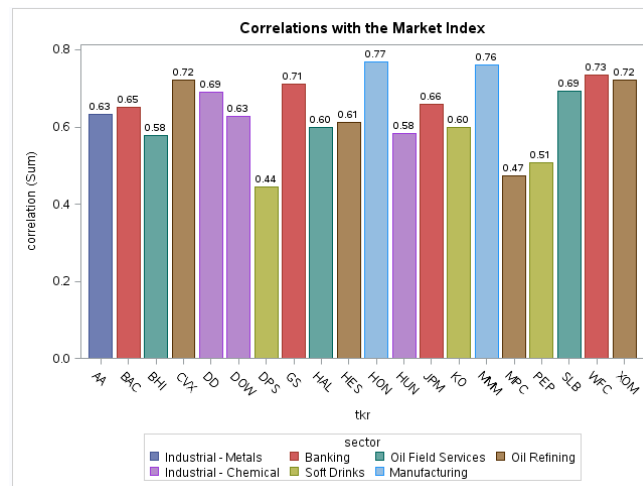
We will then perform PCA in order to deal with the multicollinearity in the data.  Next, we will create two separate regression models; the first will include all of the daily log return variables against the log return of the fund (response_VV) and the second will leverage eight components we will get from our PCA analysis.  We will then leverage a test data set to assess the predictive accuracy of each model. Lastly, we will assess the goodness-of-fit for both models and compare several metrics (adjusted R-squared, MAE, MSE) on both the training and test datasets for each model. The final goal is to determine if PCA improves the predictive accuracy of the model by removing the multicollinearity that may or may not exist among the predictor variables.

**Correlation Analysis**

The first step in our analysis is to perform the log transformation on the twenty stocks and the Vanguard market index fund.  Next, we will plot a Pearson correlation matrix for each stock against the market index (see Figure 1).  As you can see in Figure 1, many of the variables have a positive correlation with the market index. You can see several stocks in the same sector have similar correlations with the response variable (eg. CVX and XOM, HON and MMM).

Figure 1 – Correlation for each stock against the Market Index (VV) by Sector



**Principal Component Analysis (PCA)**

The next phase in this assignment is to perform Principal Component Analysis (PCA) on the predictor variables.  After running the PRINCOMP procedure in SAS, you can see the scree plot in Figure 3.  The Kaiser Rule tells us to drop all components with eigenvalues below 1.0, which in this case would suggest we should retain only three components.
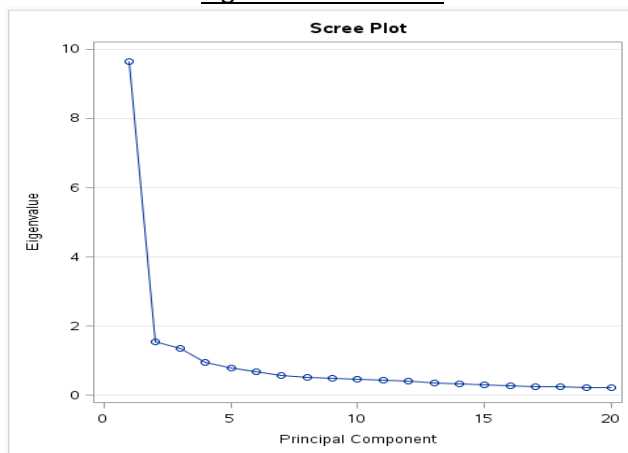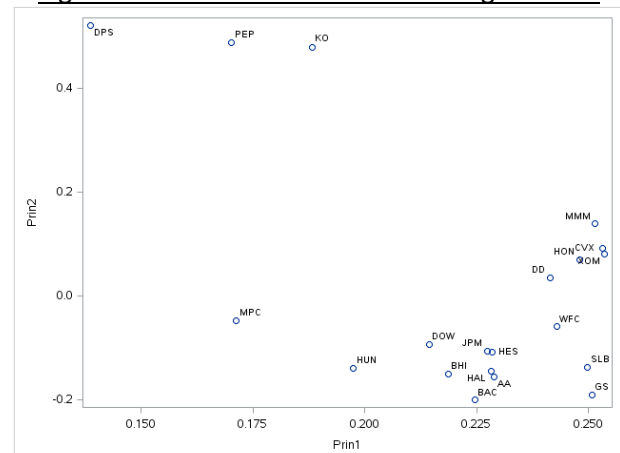
| Figure 3 – Scree Plot | Figure 4 – Scatter Plot of the first 2 Eigenvalues |
|---|---|



If you look at the chart in Figure 4, you can see that groupings have formed.  DPS, PEP, and KO all appear to fall into component 2, which makes perfect sense since they are all in the in the same sector, Soft Drinks.  You can see other similar groupings of stocks from the same industry such as CVX and XOM, both from the Oil Refining sector.

**Regression Models Leveraging Principal Component Analysis**

In this section, we will build two regression models in order to assess if performing PCA improves the predictive accuracy of the model.

*Model 1 – All Log-Returns for each Stock against the Log-Return of the Market Index - VV*

The first model we will assess include all log-return values for each stock as the predictor variables and the log-return of the Vanguard market index (response_VV).  From a goodness-of-fit perspective, the residual plots all show constant variance for all observations as shown in Figure 5.  Another indication that the model has a good fit is the straight line in the QQ plot for Model 1 (Figure 6).  Lastly, the adjusted R-squared for Model one is 0.8919, which indicates a very good fit to the data.  As can be seen in Figure 7, the variance inflation factors (VIFs) for the predictor variables are all below 5.  Thus we can conclude that no multicollinearity exists among the predictor variables.
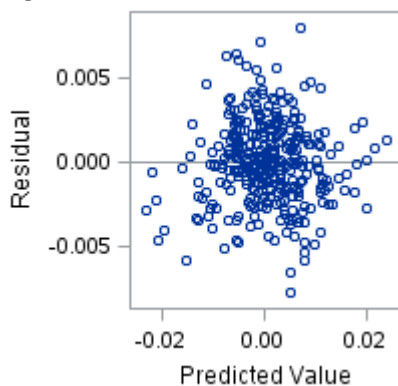
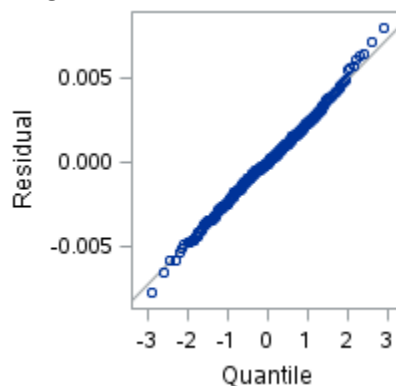Figure 5 – Residual Plot for Model 1



Figure 6 – QQ Plot for Model 1



Figure 7 – VIFs for Model 1

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| | Variance Inflation |
|---|---|---|---|---|---|---|
| Intercept | 1 | 0.00008640 | 0.00014092 | 0.61 | 0.5403 | 0 |
| return_AA | 1 | 0.01769 | 0.01317 | 1.34 | 0.1802 | 2.11490 |
| return_BAC | 1 | 0.03198 | 0.01165 | 2.75 | 0.0064 | 3.10927 |
| return_BHI | 1 | -0.00111 | 0.01323 | -0.08 | 0.9333 | 2.62997 |
| return_CVX | 1 | 0.04907 | 0.02536 | 1.93 | 0.0539 | 3.07524 |
| return_DD | 1 | 0.04674 | 0.02037 | 2.29 | 0.0224 | 2.51406 |
| return_DOW | 1 | 0.03642 | 0.01162 | 3.14 | 0.0019 | 1.88893 |
| return_DPS | 1 | 0.03670 | 0.01679 | 2.19 | 0.0295 | 1.54768 |
| return_GS | 1 | 0.04849 | 0.01555 | 3.12 | 0.0020 | 3.10450 |
| return_HAL | 1 | 0.00948 | 0.01466 | 0.65 | 0.5184 | 3.08758 |
| return_HES | 1 | 0.00359 | 0.01092 | 0.33 | 0.7425 | 2.10199 |
| return_HON | 1 | 0.12213 | 0.01924 | 6.35 | <.0001 | 2.73505 |
| return_HUN | 1 | 0.02712 | 0.00836 | 3.24 | 0.0013 | 1.79852 |
| return_JPM | 1 | 0.00902 | 0.01708 | 0.53 | 0.5979 | 3.36439 |
| return_KO | 1 | 0.07903 | 0.02226 | 3.55 | 0.0004 | 1.93633 |
| return_MMM | 1 | 0.09796 | 0.02646 | 3.70 | 0.0003 | 2.98277 |
| return_MPC | 1 | 0.01673 | 0.00809 | 2.07 | 0.0394 | 1.32999 |
| return_PEP | 1 | 0.02911 | 0.02231 | 1.30 | 0.1929 | 1.68825 |
| return_SLB | 1 | 0.03776 | 0.01709 | 2.21 | 0.0279 | 3.13690 |
| return_WFC | 1 | 0.07587 | 0.01848 | 4.10 | <.0001 | 2.59492 |
| return_XOM | 1 | 0.05467 | 0.02697 | 2.03 | 0.0435 | 2.98393 |

*Model 2 – 8 components from PCA against the Log-Return of the Market Index - VV*

In the second model, we have gone down the path of performing principal component analysis. As can be seen in Figure 8, the residuals appear to have constant variance in the model, similar to Model 1.  The QQ plot for Model 2 (Figure 9) has a straight line, which is an indication of good model fit.  Lastly, the adjusted R-squared for Model 2 is 0.8886, which is also a good indication of model fit.

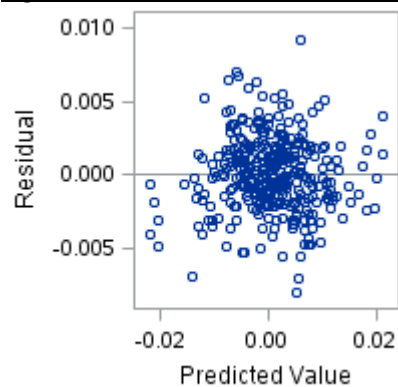Figure 8 – Residual Plot for Model 2
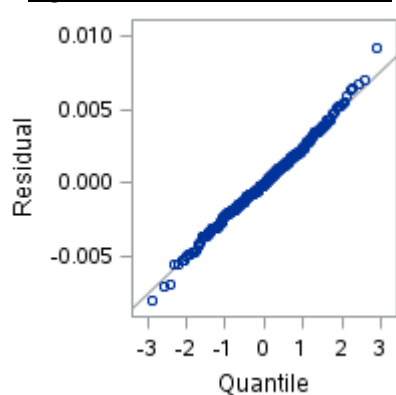


Figure 9 – QQ Plot for Model 2



Figure 10 - VIFs for Model 2

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| | Variance Inflation |
|---|---|---|---|---|---|---|
| Intercept | 1 | 0.00075978 | 0.00014045 | 5.41 | <.0001 | 0 |
| Prin1 | 1 | 0.00231 | 0.00004519 | 51.05 | <.0001 | 1.00527 |
| Prin2 | 1 | 0.00032245 | 0.00011425 | 2.82 | 0.0051 | 1.00868 |
| Prin3 | 1 | 0.00070635 | 0.00012322 | 5.73 | <.0001 | 1.00861 |
| Prin4 | 1 | 0.00030481 | 0.00014536 | 2.10 | 0.0368 | 1.00636 |
| Prin5 | 1 | -0.00017356 | 0.00015516 | -1.12 | 0.2641 | 1.00297 |
| Prin6 | 1 | 0.00000315 | 0.00017108 | 0.02 | 0.9853 | 1.00766 |
| Prin7 | 1 | -0.00010331 | 0.00018604 | -0.56 | 0.5791 | 1.02315 |
| Prin8 | 1 | -0.00040760 | 0.00020293 | -2.01 | 0.0454 | 1.02271 |

As you can see in Figure 10, the VIFs for Model 2 are all very low and do not exceed 5, indicating no multicollinearity exists among the components used in the model.

## Comparison of Model 1 and Model 2

After running the regression models on the training data, you can clearly see that both models have a good fit to the data.  However, if we review the MSE and MAE for both the training and test datasets for both models, we can see some slight differences in model performance.

<div align="center">Table 1 – Model Comparison from Training and Test samples</div>

| | | Model 1 – w/o PCA | Model 2 – w/ PCA |
|---|---|---|---|
| | Predictor(s) Selected | All stock log-returns | 8 components |
| **Training Sample** | Adjusted R2 | 0.8919 | 0.8886 |
| | MSE | 0.00000639 | 0.00000659 |
| | MAE | 0.0019020 | 0.0019752 |
| **Test Sample** | MSE | 0.00000931 | 0.00000968 |
| | MAE | 0.0021449 | 0.0021792 |

As you can see in Table 1, the adjusted R-squared, MSE, and MAE are slightly better in Model 1 than in Model 2 for both the training and test data sets.  By comparing the residual and QQ plots from both models, they almost appear to be identical.  Given we did not see any indications that multicollinearity was present among the predictor variables, this would make sense why we did not see any performance improvement from leveraging PCA.

## Predictive Accuracy & Final Model Selection

Now that we have assessed the performance of the models in the statistical sense, we will now compare the models based on their predictive accuracy.  By comparing the predicted values vs the actual values for both Model 1 and Model 2, we can have an objective basis for final model selection.

Table 2 - Model 1 Predictive Performance

| Prediction_Grade | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
| 01: Grade 1 within 10% | 20 | 12.20 | 20 | 12.20 |
| 02: Grade 2 within 10-20% | 30 | 18.29 | 50 | 30.49 |
| 03: Grade 3 within 20-30% | 18 | 10.98 | 68 | 41.46 |
| 04: Grade 4 within 30-40% | 14 | 8.54 | 82 | 50.00 |
| 05: Grade 5 above 40% | 82 | 50.00 | 164 | 100.00 |

Table 3 - Model 2 Predictive Performance

| Prediction_Grade | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
| 01: Grade 1 within 10% | 23 | 14.02 | 23 | 14.02 |
| 02: Grade 2 within 10-20% | 23 | 14.02 | 46 | 28.05 |
| 03: Grade 3 within 20-30% | 19 | 11.59 | 65 | 39.63 |
| 04: Grade 4 within 30-40% | 15 | 9.15 | 80 | 48.78 |
| 05: Grade 5 above 40% | 84 | 51.22 | 164 | 100.00 |

As we can see by comparing the frequency tables for the results of the predicted values against the actual values, there is not much of a difference between both models, with model 1 having a couple more predicted values within 40% of the actual values.  Given we did not see any multicollinearity in the data, the logical model to use would be Model 1.

## Conclusion

The final takeaway from this assignment is that there really is no benefit in performing PCA when multicollinearity does not exist among the predictor variables, at least from a predictive performance perspective. While PCA is good for dealing with multicollinearity, I can also see it being a very useful tool for identifying groupings of multiple predictors and for dimensionality reduction. The main thing I struggle with is that both of these models' predictive accuracy is far too low for me to have the confidence to recommend using either model. When we only see 50% of the predicted values coming within 40% of the actual values, I can't trust that these models are very accurate using linear regression. The next step would be to explore other machine learning techniques, such as support vector machines or naive Bayes, to see if there is any increase in predictive accuracy with the dataset we are working with in this assignment.

## SAS Code Output

```
1  *-------------------------------------------------------------------------------
2  *   Nate Bitting
3  *   Assignment 6
4  *------------------------------------------------------------------------------;
5
6  * Code used to get the data into my library;
7  ods graphics on;
8  libname mydata '/courses/d6fc9ae5ba27fe300/c_3505/SAS_Data/' access=readonly;
9  proc datasets library=mydata; run; quit;
10
11 data temp;
12 set mydata.stock_portfolio_data;
13 run;
14
15 proc sort data=temp; by date; run; quit;
16
17 data return_data;
18 set temp;
19
20 return_AA = log(AA/lag1(AA));
21 return_BAC = log(BAC/lag1(BAC));
22 return_BHI = log(BHI/lag1(BHI));
23 return_CVX = log(CVX/lag1(CVX));
24 return_DD = log(DD/lag1(DD));
25 return_DOW = log(DOW/lag1(DOW));
26 return_DPS = log(DPS/lag1(DPS));
27 return_GS = log(GS/lag1(GS));
28 return_HAL = log(HAL/lag1(HAL));
29 return_HES = log(HES/lag1(HES));
30 return_HON = log(HON/lag1(HON));
31 return_HUN = log(HUN/lag1(HUN));
32 return_JPM = log(JPM/lag1(JPM));
33 return_KO = log(KO/lag1(KO));
34 return_MMM = log(MMM/lag1(MMM));
35 return_MPC = log(MPC/lag1(MPC));
36 return_PEP = log(PEP/lag1(PEP));
37 return_SLB = log(SLB/lag1(SLB));
38 return_WFC = log(WFC/lag1(WFC));
39 return_XOM = log(XOM/lag1(XOM));
40 response_VV = log(VV/lag1(VV));
41 run;
42
43 *create a list of all the predictor variables;
44 %let xlist =return_AA return_BAC return_BHI return_CVX return_DD return_DOW return_DPS return_GS
45                  return_HAL return_HES return_HON return_HUN return_JPM return_KO
46                  return_MMM return_MPC return_PEP return_SLB return_WFC return_XOM;
47
48 proc print data=return_data(obs=10); run; quit;
49
50 ods output PearsonCorr=portfolio_correlations;
51 proc corr data=return_data;
52 var return_:;
53 with response_VV;
54 run; quit;
55
56 proc print data=portfolio_correlations; run; quit;
57
58 data wide_correlations;
59 set portfolio_correlations (keep=return_:);
60 run;
61
62 proc transpose data=wide_correlations out=long_correlations;
63 run; quit;
64
65 data long_correlations;
66 set long_correlations;
67 tkr = substr(_NAME_,8,3);
68 drop _NAME_;
69 rename COL1=correlation;
70 run;
71
72 proc print data=long_correlations; run; quit;
73
74 *print a scatter plot of a few variables with high correlation with the response;
75 proc sgscatter data=return_data;
76 title 'Scatter Plots of a few Predictors';
77 plot return_HON*return_MMM return_GS*return_XOM;
78 run;
```

```
80  data sector;
81  input tkr $ 1-3 sector $ 4-35;
82  datalines;
83  AA Industrial - Metals
84  BAC Banking
85  BHI Oil Field Services
86  CVX Oil Refining
87  DD Industrial - Chemical
88  DOW Industrial - Chemical
89  DPS Soft Drinks
90  GS Banking
91  HAL Oil Field Services
92  HES Oil Refining
93  HON Manufacturing
94  HUN Industrial - Chemical
95  JPM Banking
96  KO Soft Drinks
97  MMM Manufacturing
98  MPC Oil Refining
99  PEP Soft Drinks
100 SLB Oil Field Services
101 WFC Banking
102 XOM Oil Refining
103 VV Market Index
104 ;
105 run;
106
107 proc print data=sector; run; quit;
108
109 proc sort data=sector; by tkr; run;
110
111 proc sort data=long_correlations; by tkr; run;
112
113 data long_correlations;
114 merge long_correlations (in=a) sector (in=b);
115 by tkr;
116 if (a=1) and (b=1);
117 run;
118
119 proc print data=long_correlations; run; quit;
120
121 ods graphics on;
122 title 'Correlations with the Market Index';
123 proc sgplot data=long_correlations;
124 format correlation 3.2;
125 vbar tkr / response=correlation group=sector groupdisplay=cluster
126 datalabel;
127 run; quit;
128 ods graphics off;
129
130
131 * create a single dataset that only includes the log return values for the predictor and repsonse variable;
132 data return_data_only;
133 set return_data;
134 drop AA;
135 drop BAC;
136 drop BHI;
137 drop CVX;
138 drop DD;
139 drop DOW;
140 drop DPS;
141 drop GS;
142 drop HAL;
143 drop HES;
144 drop HON
145 drop HUN;
146 drop JPM;
147 drop KO;
148 drop MMM;
149 drop MPC;
150 drop PEP;
151 drop SLB;
152 drop WFC;
153 drop XOM;
154 drop VV;
155 run;
```

```sas
157  ods graphics on;
158  proc princomp
159  data=return_data_only
160  out=pca_output
161  outstat=eigenvectors
162  plots=scree(unpackpanel);
163  var &xlist;
164  run; quit;
165  ods graphics off;
166
167  proc print data=pca_output(obs=10); run;
168
169  proc print data=eigenvectors(where=(_TYPE_='SCORE')); run;
170
171  data pca2;
172  set eigenvectors(where=(_NAME_ in ('Prin1','Prin2')));
173  drop _TYPE_ ;
174  run;
175
176  proc print data=pca2; run;
177
178  proc transpose data=pca2 out=long_pca; run; quit;
179  proc print data=long_pca; run;
180
181  data long_pca;
182  set long_pca;
183  format tkr $3.;
184  tkr = substr(_NAME_,8,3);
185  drop _NAME_;
186  run;
187
188  proc print data=long_pca; run;
189
190  * Plot the first two principal components;
191  ods graphics on;
192  proc sgplot data=long_pca;
193  scatter x=Prin1 y=Prin2 / datalabel=tkr;
194  run; quit;
195  ods graphics off;
196
197  *******************************************************************;
198  * Create a training data set and a testing data set from the
199  PCA output;
200  * Note that we will use a SAS shortcut to keep both of these
201  'datasets' in one data set that we will call cv_data (cross-validation
202  data). ;
203  *******************************************************************;
204  data cv_data;
205  merge pca_output return_data_only(keep=response_VV);
206  * No BY statement needed here. We are going to append a column in
207  its current order;
208  * generate a uniform(0,1) random variable with seed set to
209  123; u = uniform(123);
210  if (u < 0.70) then train = 1;
211  else train = 0;
212  if (train=1) then train_response=response_VV;
213  else train_response=.;
214  run;
215  proc print data=cv_data(obs=10); run;
216
217
218  *regression model 1 without PCA;
219  proc reg data=cv_data outest=RegOut;
220  model train_response = &xlist / vif mse;
221  *output out=residuals_final (keep = resid_final) r=resid_final;
222  run;
223
224  * calc the MAE for Model 1;
225  data abs_resid;
226      set residuals_final;
227      abs_resid = abs(resid_final);
228  run;
229
230  proc means data=abs_resid mean;
231      var abs_resid;
232  run;
233
234  * calcualte the estimated values using the test dataset;
```

```sas
235 proc score data=cv_data score=RegOut out=RScoreP type=parms;
236     var &xlist;
237 run;
238
239 proc score data=cv_data score=RegOut out=RScoreR type=parms;
240     var train_response &xlist;
241 run;
242
243 * Output the scores using the model built with the training dataset against the test dataset;
244 proc score data=cv_data score=RegOut out=NewPred_M1 type=parms
245             nostd predict;
246     var train_response &xlist;
247 run;
248
249 * Smart data formats for predictions;
250 proc format;
251 value pred_acc_sfmt
252     0 -< .1 = '01: Grade 1 within 10%'
253    .1 -< .2 = '02: Grade 2 within 10-20%'
254    .2 -< .3 = '03: Grade 3 within 20-30%'
255    .3 -< .4 = '04: Grade 4 within 30-40%'
256    other = '05: Grade 5 above 40%'
257    ;
258 run;
259
260 *create a new dataset that contains the test set data and the predictive scores for Model 1;
261 data prediction_Data_M1;
262     set NewPred_M1;
263     if (train_response = null);
264     pred_score = abs(response_VV / MODEL1 - 1);
265     abs_resid = abs(response_VV - MODEL1);
266     error_term = response_VV - MODEL1;
267     sq_error = error_term**2;
268     Prediction_Grade = put(pred_score,pred_acc_sfmt.);
269 run;
270
271 * calculate the MAE and MSE from the test sample;
272 proc means data=prediction_Data_M1 mean;
273     var abs_resid sq_error;
274 run;
275
276 * create a frequency table to show the operational accuracy of model 1;
277 proc freq data=prediction_Data_M1;
278     TABLES Prediction_Grade;
279 run;
280
281
282 *regression model 2 with 8 components;
283 proc reg data=cv_data outest=RegOutPCA1;
284 model train_response = Prin1 Prin2 Prin3 Prin4 Prin5 Prin6 Prin7 Prin8 / vif mse;
285 *output out=residuals_PCA1 (keep = resid_PCA1) r=resid_PCA1;
286 run;
287
288 * calc the MAE for Model 2;
289 data abs_resid_PCA1;
290     set residuals_PCA1 ;
291     abs_resid_PCA1 = abs(resid_PCA1);
292 run;
293
294 proc means data=abs_resid_PCA1 mean;
295     var abs_resid_PCA1;
296 run;
297
298 * calcualte the estimated values using the test dataset;
299 proc score data=cv_data score=RegOutPCA1 out=RScoreP_M2 type=parms;
300     var Prin1 Prin2 Prin3 Prin4 Prin5 Prin6 Prin7 Prin8;
301 run;
302
303 proc score data=cv_data score=RegOutPCA1 out=RScoreR_M2 type=parms;
304     var train_response Prin1 Prin2 Prin3 Prin4 Prin5 Prin6 Prin7 Prin8;
305 run;
306
307 * Output the scores using the model built with the training dataset against the test dataset;
308 proc score data=cv_data score=RegOutPCA1 out=NewPred_M2 type=parms
309             nostd predict;
310     var train_response Prin1 Prin2 Prin3 Prin4 Prin5 Prin6 Prin7 Prin8;
311 run;
```

```sas
314  *create a new dataset that contains the test set data and the predictive scores for Model 2;
315  data prediction_Data_M2;
316       set NewPred_M2;
317       if (train_response = null);
318       pred_score = abs(response_VV / MODEL1 - 1);
319       abs_resid = abs(response_VV - MODEL1);
320       error_term = response_VV - MODEL1;
321       sq_error = error_term**2;
322       Prediction_Grade = put(pred_score,pred_acc_sfmt.);
323  run;
324
325  * calculate the MAE and MSE from the test sample;
326  proc means data=prediction_Data_M2 mean;
327       var abs_resid sq_error;
328  run;
329
330  * create a frequency table to show the operational accuracy of model 2;
331  proc freq data=prediction_Data_M2;
332       TABLES Prediction_Grade;
333  run;
334
335  *regression model 3 with 3 components;
336  proc reg data=cv_data outest=RegOutPCA2;
337  model train_response = Prin1 Prin2 Prin3 / vif mse;
338  output out=residuals_PCA2 (keep = resid_PCA2) r=resid_PCA2;
339  run;
340
341  * calc the MAE for Model 3;
342  data abs_resid_PCA2;
343       set residuals_PCA2 ;
344       abs_resid_PCA2 = abs(resid_PCA2);
345  run;
346
347  proc means data=abs_resid_PCA2 mean;
348       var abs_resid_PCA2;
349  run;
350
351  * calcualte the estimated values using the test dataset;
352  proc score data=cv_data score=RegOutPCA2 out=RScoreP_M3 type=parms;
353       var Prin1 Prin2 Prin3;
354  run;
355
356  proc score data=cv_data score=RegOutPCA2 out=RScoreR_M3 type=parms;
357       var train_response Prin1 Prin2 Prin3;
358  run;
359
360  * Output the scores using the model built with the training dataset against the test dataset;
361  proc score data=cv_data score=RegOutPCA2 out=NewPred_M3 type=parms
362               nostd predict;
363       var train_response Prin1 Prin2 Prin3;
364  run;
365
366
367  *create a new dataset that contains the test set data and the predictive scores for Model 3;
368  data prediction_Data_M3;
369       set NewPred_M3;
370       if (train_response = null);
371       pred_score = abs(response_VV / MODEL1 - 1);
372       abs_resid = abs(response_VV - MODEL1);
373       error_term = response_VV - MODEL1;
374       sq_error = error_term**2;
375       Prediction_Grade = put(pred_score,pred_acc_sfmt.);
376  run;
377
378  * calculate the MAE and MSE from the test sample;
379  proc means data=prediction_Data_M3 mean;
380       var abs_resid sq_error;
381  run;
382
383  * create a frequency table to show the operational accuracy of model 3;
384  proc freq data=prediction_Data_M3;
385       TABLES Prediction_Grade;
386  run;
```