

```

1 *****
2 * TITLE : SAS GRAIN PRICE PROJECT ANALYSIS
3 *
4 * DESCRIPTION: Final project for BIOS 7400 with Xiao Song, UGA, Spring 2022.
5 * Simple analysis of grain price data.
6 *
7 *-----
8 * JOB NAME: analysis.SAS
9 * LANGUAGE: SAS v9.4 (on demand for academics)
10 *
11 * NAME: Zane Billings
12 * DATE: 2022-04-22
13 *
14 *****;
15
16 FOOTNOTE "Job run by Zane Billings on &SYSDATE at &SYSTIME.";
17
18 TITLE 'ANALYSIS OF USDA HISTORICAL GRAIN PRICE DATA';
19
20 OPTIONS NODATE LS=95 PS=42;
21
22 LIBNAME HOME '/home/u59465388/SAS-Grain-Prices';
23
24 ODS GRAPHICS / WIDTH = 6in HEIGHT = 3in;
25
26 *****;
27 * Show the descriptor portion of the dataset;
28 *****;
29
30 TITLE2 "CONTENTS OF GRAINS DATASET";
31
32 -----
33 PROC CONTENTS DATA = HOME.GRAINS;
34 RUN;
35
36 *****;
37 * Plot outcome time series;
38 *****;
39
40 FOOTNOTE; * Remove the footnote so it isn't on the graphs;
41
42 * Plot the time series of log grain price over time. This makes a separate
43 time series line for each grain.;
44 TITLE2 "PRICE PER BUSHEL OF GRAINS OVER TIME";
45 -----
46 PROC SGPLOT DATA = HOME.GRAINS;
47 SERIES X = YEAR Y = LPE / GROUP = GRN;
48 RUN;
49
50 * Color the points of the time series by the President's political party. The
51 default colors are already red and blue so we don't need to change them!
52 Also plots a gray line underneath the points, since the JOIN option for
53 the SCATTER statement will not connect the points in order.;
54 TITLE2 "GRAIN PRICES AND PRESIDENT'S POLITICAL PARTY OVER TIME";
55 -----
56 PROC SGPANEL DATA = HOME.GRAINS;
57 PANELBY GRN;

```

```

56     SERIES X = YEAR Y = LPE / LINEATTRS = (COLOR = "GRAY") SMOOTHCONNECT;
57     SCATTER X = YEAR Y = LPE / GROUP = PARTY
58     MARKERATTRS = (SYMBOL = CIRCLEFILLED);
59 RUN;
60
61 ODS GRAPHICS / WIDTH = 6in HEIGHT = 6in;
62
63 * Make a boxplot of log price vs. president's political party. This ignores the
64   time series information, but can tell us if either party has more high or
65   low years compared to the other.;
66 TITLE2 "LOG PRICE DISTRIBUTION BY PRESIDENT'S POLITICAL PARTY";
67 PROC SGPNL DATA = HOME.GRAINS;
68     PANELBY GRN;
69     HBOX LPE / GROUP = PARTY;
70 RUN;
71
72 ODS GRAPHICS / WIDTH = 6in HEIGHT = 9in;
73
74 * Make a scatterplot of the log price vs each covariate, ignoring the time
75   series component of the data. There is not an easy way to connect the points
76   like a phase portrait using PROC SGPLOT.;
77 * I divided this into two plots so they would fit on one page nicer. In the final
78   manuscript they could be put side by side.;
79 TITLE2 "SCATTERPLOTS OF PRICE VS COVARIATES";
80 PROC SGSCATTER DATA = HOME.GRAINS;
81     PLOT LPE * (ACR HVT LNR PRD YLD) / REG
82     COLUMNS = 2 GROUP = GRN;
83 RUN;
84
85 PROC SGSCATTER DATA = HOME.GRAINS;
86     PLOT LPE * (INFL PWR TEMP VALUE) / REG
87     COLUMNS = 2;
88 RUN;
89
90 *****;
91 * Plots of covariates across time;
92 *****;
93
94 * Plot the time series of each covariate, to assess how they change. I split
95   this one into two plots to prevent the plots being too small as before.;
96 TITLE2 "CHANGE IN COVARIATES ACROSS TIME";
97 PROC SGSCATTER DATA = HOME.GRAINS;
98     PLOT (ACR HVT LNR PRD YLD) * YEAR /
99     COLUMNS = 2 GROUP = GRN JOIN MARKERATTRS = (SIZE = 0);
100 RUN;
101
102
103 PROC SGSCATTER DATA = HOME.GRAINS;
104     PLOT (INFL PWR TEMP VALUE) * YEAR /
105     COLUMNS = 2 JOIN MARKERATTRS = (SIZE = 0);
106 RUN;
107
108 *****;
109 * Univariate analyses;
110 *****;
111

```

```

112 * Univariate analysis of main outcome (log price) by grain type;
113 TITLE2 "UNIVARIATE SUMMARY OF GRAIN DATA OVER TIME";
114
115 ODS GRAPHICS / WIDTH = 4in HEIGHT = 4in;
116
117 PROC UNIVARIATE DATA = HOME.GRAINS PLOTS;
118     VAR LPE;
119     CLASS GRN;
120 RUN;
121
122 *****;
123 * Bivariate analyses of price and covariates, ignoring time;
124 *****;
125
126 TITLE2 "BIVARIATE CORRELATIONS ACROSS NUMERICAL VARIABLES";
127 * Correlations -- check to see which covariates are correlated with the outcome,
128     and which are correlated with each other and should not be modeled
129     together.;
130 PROC CORR PEARSON SPEARMAN DATA = HOME.GRAINS;
131     VAR LPE ACR HVT LNR PRD YLD INFL PWR TEMP VALUE;
132     BY GRN;
133 RUN;
134
135 TITLE2 "SUMMARY STATISTICS BY PRESIDENTIAL PARTY AND GRAIN";
136 * Mean difference in LPE by party -- proc corr does not have a point biserial
137     option, so we can check the difference/overlap in means and standard errors
138     to assess if party seems to impact log price for any of the grains.;
139 PROC MEANS DATA = HOME.GRAINS MEAN STDERR MEDIAN RANGE NWAY;
140     VAR LPE;
141     CLASS PARTY;
142     BY GRN;
143 RUN;
144
145 *****;
146 * Simple and multiple OLS regression models;
147 *****;
148
149 TITLE2 "SIMPLE LINEAR REGRESSION MODELS";
150
151 * Model stratified by grain only;
152 PROC GLM DATA = HOME.GRAINS PLOTS = ALL;
153     CLASS GRN;
154     MODEL LPE = GRN / NOINT;
155 RUN;
156
157 * Write a macro to fit all regression models of the form
158     MODEL COVAR GRN COVAR * GRN
159     without having to type out all of the PROC GLM statements. This model will
160     be parametrized without an intercept, and will generate all appropriate
161     diagnostic plots for the model.;
162
163
164 %MACRO ALLSIMPLE(DAT = , RESP = , PRED = );
165     %LET N = %SYSFUNC(COUNTW(&PRED));
166     %DO I = 1 %TO &N;
167         PROC GLM DATA = &DAT PLOTS = ALL;

```

```

168         CLASS GRN;
169         MODEL &RESP = %SCAN(&PRED, &I) | GRN / NOINT;
170     RUN;
171 %END;
172 %MEND;
173
174 %ALLSIMPLE(
175     DAT = HOME.GRAINS,
176     RESP = LPE,
177     PRED = ACR PRD INFL TEMP PWR YEAR
178 );
179
180 * Fit the same model that was used as before, but with party as a covariate.
181   Party needs to be in the class statement, and is the only categorical
182   variable, so it wasn't worth modifying the above macro to use party
183   correctly and I did it manually.;
184 PROC GLM DATA = HOME.GRAINS PLOTS = ALL;
185     CLASS GRN PARTY;
186     MODEL LPE = GRN | PARTY / NOINT;
187 RUN;
188
189 TITLE2 "1866 FULL MODEL";
190 * 1866 FULL MODEL: this model includes all non-correlated predictors that were
191   measured in 1866.;
192 PROC MIXED DATA = HOME.GRAINS PLOTS = ALL;
193     CLASS GRN PARTY;
194     MODEL LPE =
195         GRN HVT PRD INFL PWR YEAR PARTY
196         GRN*HVT GRN*PRD GRN*INFL GRN*PWR GRN*YEAR GRN*PARTY /
197         NOINT SOLUTION
198     ;
199 RUN;
200
201 TITLE2 "1880 FULL MODEL";
202 * FULL MODEL WITH TEMP (1880 MODEL): this model is the same as the previous
203   model, but also includes the temperature anomaly information. Consequently,
204   it only uses data from 1880 onwards (even less for sorghum).;
205 PROC MIXED DATA = HOME.GRAINS PLOTS = ALL;
206     CLASS PARTY GRN;
207     MODEL LPE =
208         GRN HVT PRD INFL PWR YEAR PARTY TEMP
209         GRN*HVT GRN*PRD GRN*INFL GRN*PWR GRN*YEAR GRN*PARTY TEMP*PARTY /
210         NOINT SOLUTION
211     ;
212 RUN;
213
214 *****;
215 * GLS multiple regression analysis;
216 *****;
217
218 * Take the better fitting (by AIC) of the two previous models, and run a model
219   that can account for correlation using generalized least squares.
220   This model assumes exchangeable correlations between each of the time points.;
221 TITLE2 "GENERALIZED LEAST SQUARES MODEL";
222 PROC MIXED DATA = HOME.GRAINS PLOTS = ALL;
223

```

```

224 CLASS GRN;
225 MODEL LPE = HVT PRD INFL PWR YEAR GRN GRN*HVT GRN*PRD /
226 NOINT SOLUTION CHISQ;
227 REPEATED;
228 RUN;
229
230 *****;
231 * Simple forecasting;
232 *****;
233
234 * Now instead of just using regression models, we can try to fit a more
235 flexible forecasting model using PROC ARIMA.
236 First, we need a time variable that is actually a SAS date, so we create
237 that first.;
238 DATA TS_DAT;
239 SET HOME.GRAINS;
240 T = MDY(1, 1, YEAR);
241 RUN;
242
243 * Next we use the IDENTIFY modeling stage. We check up to 30 lags in the
244 first ARIMA modeling stage, and also explicitly test for stationarity at
245 the first 10 differences using the random walk with drift test. We
246 also use the SCAN method, which is a heuristic for identifying
247 candidate ARIMA models.;
248 PROC ARIMA DATA = TS_DAT;
249 IDENTIFY VAR = LPE NLAG = 30 SCAN STATIONARITY = (RW = 10);
250 BY GRN;
251 TITLE2 "ARIMA TESTS";
252 RUN;
253
254 * Next we use the ESTIMATE modeling stage. We fit several different ARIMA
255 models to the data in order to see which fits our time series the best,
256 and if any have white noise as the error term.;
257 * One PROC ARIMA can contain multiple ESTIMATE statements, but I split these
258 into multiple PROC steps to make the output easier to read.;
259 * We are basically fitting all of these models to get the AIC and see which is
260 the best fit.;
261
262 * Model 1: AR(1);
263 PROC ARIMA DATA = TS_DAT;
264 IDENTIFY VAR = LPE;
265 ESTIMATE P = 1;
266 BY GRN;
267 TITLE2 "AR(1)";
268 RUN;
269
270 * MODEL 2: AR(2);
271 PROC ARIMA DATA = TS_DAT;
272 IDENTIFY VAR = LPE;
273 ESTIMATE P = 2;
274 BY GRN;
275 TITLE2 "AR(2)";
276 RUN;
277
278 * MODEL 3: MA(1);
279

```

```
280 PROC ARIMA DATA = TS_DAT;  
281     IDENTIFY VAR = LPE;  
282     ESTIMATE Q = 1;  
283     BY GRN;  
284     TITLE2 "MA(1)";  
285 RUN;  
286  
287 * MODEL 4: ARMA(1, 1);  
288 PROC ARIMA DATA = TS_DAT;  
289     IDENTIFY VAR = LPE;  
290     ESTIMATE P = 1 Q = 1;  
291     BY GRN;  
292     TITLE2 "ARMA(1, 1)";  
293 RUN;  
294  
295 * MODEL 5: ARIMA(1, 1, 0);  
296 PROC ARIMA DATA = TS_DAT;  
297     IDENTIFY VAR = LPE(1);  
298     ESTIMATE P = 1;  
299     BY GRN;  
300     TITLE2 "ARIMA(1, 1, 0)";  
301 RUN;  
302  
303 * MODEL 6: ARIMA(1, 1, 1);  
304 PROC ARIMA DATA = TS_DAT;  
305     IDENTIFY VAR = LPE(1);  
306     ESTIMATE P = 1 Q = 1;  
307     BY GRN;  
308     TITLE2 "ARIMA(1, 1, 1)";  
309 RUN;  
310  
311 * MODEL 7: ARIMA(0,0,0) (WHITE NOISE);  
312 PROC ARIMA DATA = TS_DAT;  
313     IDENTIFY VAR = LPE;  
314     ESTIMATE P = 0 Q = 0;  
315     BY GRN;  
316     TITLE2 "ARIMA(0, 0, 0)";  
317 RUN;  
318  
319 * MODEL 8: ARIMA(0,1,0) (RANDOM WALK);  
320 PROC ARIMA DATA = TS_DAT;  
321     IDENTIFY VAR = LPE(1);  
322     ESTIMATE P = 0 Q = 0;  
323     BY GRN;  
324     TITLE2 "ARIMA(0, 1, 0)";  
325 RUN;  
326  
327  
328 * Finally, we use the best fitting model to make some simple forecasts in  
329   the FORECAST modeling stage. We also identify outliers of the best  
330   fitting model.;  
331  
332 PROC ARIMA DATA = TS_DAT;  
333     IDENTIFY VAR = LPE(1);  
334     ESTIMATE P = 1 Q = 1;  
335     OUTLIER;
```

```
336     FORECAST LEAD = 10 INTERVAL = YEAR ID = T OUT = GRAIN_FC;  
337     BY GRN;  
338     TITLE2 "FORECASTING";  
339 RUN;  
340  
341  
342  
343  
...
```