```
1   ************************************************************************
2   *  TITLE :      SAS GRAIN PRICE PROJECT
3   *
4   *  DESCRIPTION: Final project for BIOS 7400 with Xiao Song, UGA, Spring 2022.
5   *               Cleaning data for grain price analysis.
6   *
7   *---------------------------------------------------------------------
8   *  JOB NAME:    cleaning.SAS
9   *  LANGUAGE:    SAS v9.4 (on demand for academics)
10  *
11  *  NAME:        Zane Billings
12  *  DATE:        2022-04-20
13  *
14  ************************************************************************;
15
16  FOOTNOTE "Job run by Zane Billings on &SYSDATE at &SYSTIME";
17
18  TITLE 'Grain Price Analysis';
19
20  OPTIONS NODATE LS=95 PS=42;
21
22  LIBNAME HOME '/home/u59465388/SAS-Grain-Prices';
23
24  ************************************************************************;
25  * Macros;
26  ************************************************************************;
27
28  * Variables for filtering the years to export in the cleaned dataset. I have
29      them set to the min/max values in the dataset, but this allows for easier
30      changing than specifying the years manually.;
31  %LET MINYEAR = 1866;
32  %LET MAXYEAR = 2021;
33
34  * Variable for controlling whether the following macro prints to the report.
35      It is easier to toggle this in one place than to add or remove the macro
36      calls later in the script.
37      1: Prints first &PRINTN observations of the dataset and the descriptor
38          portion as well.
39      Any other value (preferably 0): does not print (indeed, the macro will
40          not execute anything after the logical step).;
41  %LET VERBOSE = 1;
42  %LET PRINTN = 10;
43
44  * Macro for printing values and descriptor portion of data;
45  %MACRO DESCRIBE (DAT =, N = &PRINTN);
46      %IF %EVAL(&VERBOSE = 1) %THEN %DO;
47          PROC PRINT DATA = &DAT (OBS = &N) LABEL;
48          RUN;
49
50          PROC CONTENTS DATA = &DAT;
51          RUN;
52      %END;
53  %MEND;
54
55  ************************************************************************;
56  * Data importing;
57
```

```sas
58  **************************************************************************;
59
60  * Import the temperature anomaly data;
61  FILENAME NASATEMP "/home/u59465388/SAS-Grain-Prices/nasatemp.txt";
62  DATA TEMP;
63      * Read in the NASA temperature data. The data starts at line 9.;
64      INFILE NASATEMP FIRSTOBS = 9;
65
66      * Bring the next line of the INFILE into the input buffer;
67      INPUT @;
68
69      * If the first detectable word (which should be the YEAR) is not a numeric
70        digit, delete the row from the buffer, and thus do not import it.
71        This skips the blank rows and repeated header rows.
72        After DELETE is executed, return to the beginning of the data step.;
73      IF NOTDIGIT(SCAN(_INFILE_, 1)) THEN DELETE;
74
75      * If the YEAR is a number, import the current infile into the dataset;
76      ELSE DO;
77          * The data has missing values coded as '****', replace these with . so that
78            SAS interprets them as missing correctly.;
79          _INFILE_ = TRANSRN(_INFILE_, "****", ".");
80          * Read in only the first 13 columns.;
81          INPUT YEAR JAN FEB MAR APR MAY JUN JUL AUG SEP OCT NOV DEC;
82      END;
83
84      * Get the yearly average, and then divide by 100 to make the units degrees C.
85        Round to two decimal places.;
86      TEMP = ROUND(MEAN(OF JAN -- DEC) / 100, 0.01);
87      DROP JAN -- DEC;
88
89      * Give information labels to the variables;
90      LABEL
91          YEAR = "Calendar year"
92          TEMP = "Temperature diff. (deg. C)"
93      ;
94  RUN;
95
96  %DESCRIBE(DAT = WORK.TEMP);
97
98  * Import the presidential party data;
99  FILENAME PRESI '/home/u59465388/SAS-Grain-Prices/presidential.csv';
100 DATA PRES;
101     * Set length of variables to ensure character vars don't get cut off;
102     LENGTH YEAR 4 PRES $ 20 PARTY $ 25;
103
104     * Import CSV file, nothing complicated like the last file;
105     INFILE PRESI DLM = ',' FIRSTOBS = 2;
106     INPUT YEAR PRES $ PARTY $;
107
108     * Abraham Lincoln and Andrew Johnson are listed as 'National Union' party
109       members, but this isn't terribly useful. Historically, Abraham Lincoln
110       was a Republican and Andrew Johnson was a Democrat, and the National Union
111       coalition was a transitional step. So I'll recode these two for simplicity.;
112     IF PRES = "Abraham Lincoln" THEN PARTY = "Republican";
113     ELSE IF PRES = "Andrew Johnson" THEN PARTY = "Democrat";
114
115
```

```sas
116        * Add descriptive labels;
117        LABEL
118            YEAR = "Calendar year"
119            PRES = "President name"
120            PARTY = "President party"
121        ;
122    RUN;
123
124    * The presidential data only goes through 2013, so we will have to manually
125        input the 2013 - 2022 data and append that to the end.;
126    DATA PRES_END;
127        LENGTH YEAR 4 PRES $ 20 PARTY $ 25;
128        INPUT YEAR PRES $ PARTY $;
129        LABEL
130            YEAR = "Calendar year"
131            PRES = "President name"
132            PARTY = "President party"
133        ;
134        INFILE DATALINES DSD DLM = " ";
135        DATALINES;
136    2014 "Barack Obama" "Democrat"
137    2015 "Barack Obama" "Democrat"
138    2016 "Barack Obama" "Democrat"
139    2017 "Donald Trump" "Republican"
140    2018 "Donald Trump" "Republican"
141    2019 "Donald Trump" "Republican"
142    2020 "Donald Trump" "Republican"
143    2021 "Joseph Biden" "Democrat"
144    2022 "Joseph Biden" "Democrat"
145    ;
146    RUN;
147
148    * Now append the second dataset to the end of the first;
149    PROC APPEND BASE = WORK.PRES DATA = WORK.PRES_END;
150    RUN;
151
152    %DESCRIBE(DAT = WORK.PRES);
153
154    * Import the inflation data;
155    FILENAME INFL '/home/u59465388/SAS-Grain-Prices/inflation_data.csv';
156    DATA INFLATION;
157        * Import CSV file, easy like the presidential data;
158        INFILE INFL DLM = ',' FIRSTOBS = 2;
159        INPUT YEAR VALUE INFL;
160
161        * Create a new column for relative 'worth': 1 / value in 1886 dollars
162          is the 'buying power' of $1 relative to an 1866 dollar.;
163        PWR = ROUND(1 / VALUE, 0.01);
164
165        * Assign descriptive lables;
166        LABEL
167            YEAR = 'Calendar year'
168            VALUE = 'Adjusted value'
169            INFL = 'Rate of inflation'
170            PWR = 'Buying power'
171        ;
172    RUN;
173
```

```sas
174
175  %DESCRIBE(DAT = WORK.INFLATION);
176
177  * Import the feed grains data. This is a complex and messy excel spreadsheet
178      that is easy to manually view but difficult to use as actual data. For
179      this project, I will only clean the first sheet.;
180  * In the current form, importing the data will be quite complicated and I think
181      impossible using PROC IMPORT. So I opened the dataset in Excel and exported
182      the sheet that I needed as a CSV file, which is what I'll import here.;
183  FILENAME FDGRN '/home/u59465388/SAS-Grain-Prices/fg-sheet1.csv';
184
185  DATA ALLGRNS;
186      * Import the CSV file. The option DSD is necessary to read in consecutive
187          delimiters as missing data, and the MISSOVER option is necessary as
188          there are missing values at the end of lines, so the INPUT specification
189          should be interpreted strictly.;
190      INFILE FDGRN DLM = ',' FIRSTOBS = 9 DSD MISSOVER;
191
192      * SAS doesn't like the missing values being denoted by ,, even with the DSD
193          option, and has a hard time parsing the numeric values. So, I'll import
194          all of the variables as character variables with silly names. The
195          names are uninformative, but easy to use all together in SAS statements.
196        Note that I have also included the trailing @ so I can check the next line
197          for all blanks, and delete the line before being read if that is the case.;
198      INPUT GRN $ YR $ V1 $ V2 $ V3 $ V4 $ V5 $ V6 $ @;
199
200      * If the next line (@) is all missing, do not read it in;
201      IF MISSING(YR) THEN DELETE;
202
203      * The grain variable is only denoted once, and is missing for all other
204          records in the time series. This part of the code saves the most recent
205          non-missing value of GRN, and then uses it to fill in the value of
206          all missing GRN values until it finds a new non-missing value.;
207      IF NOT MISSING(GRN) THEN DO;
208          TMP = GRN;
209          RETAIN TMP;
210      END;
211      ELSE GRN = TMP;
212
213      * Create a YEAR variable as the first four digits of the YR variable, which
214          looks like ####/##. Use INPUT() to make this new variable numeric.;
215      YEAR = INPUT(SUBSTR(YR, 1, 4), 4.);
216
217      * Convert the imported character variables to numeric variables. Since SAS
218          cannot modify variable types in place, we have to create two arrays. One
219          array (_CHA) holds the placeholder character variables, and the second array
220          (_NUM) holds the newly declared numeric variables with somewhat better
221          names. Then we handle the missing character values explicitly to prevent SAS
222          from complaining about the blanks, and use INPUT to parse the remaining
223          values to numbers. We use the comma informat here since some of the
224          numeric values have commas as place value separators.;
225      ARRAY _CHA{6} $ V1 - V6;
226      ARRAY _NUM{6} ACR HVT PRD YLD PCE LNR;
227      DO I = 1 TO 6;
228          IF MISSING(_CHA{I}) THEN _NUM{I} = .;
229          ELSE _NUM{I} = INPUT(_CHA{I}, COMMA8.);
230      END;
231
```

```sas
232
233         * Compute the percent change from the previous year;
234         PCT = ROUND(DIF(PCE) / LAG(PCE) * 100, 0.01);
235
236         * Compute the log of the price;
237         LPE = LOG10(PCE);
238
239         * Drop all of the temporary and placeholder variables that we don't need in
240             the cleaned dataset;
241         DROP TMP YR V1 - V6 I;
242
243         * Assign descriptive labels to the remaining useful variables.;
244         LABEL
245             GRN = "Grain commodity"
246             YEAR = "Calendar year"
247             ACR = "Acerage (M)"
248             HVT = "Acres harvested (M)"
249             PRD = "Bushels produced (M)"
250             YLD = "Yield (bushels per acre)"
251             PCE = "Price per bushel"
252             LPE = "log10 price per bushel"
253             LNR = "Loan rate per bushel"
254             PCT = "Pct change in price"
255         ;
256 RUN;
257
258 %DESCRIBE(DAT = WORK.ALLGRNS);
259
260 *****************************************************************************;
261 * Data merging;
262 *****************************************************************************;
263
264
265 * Next, we need to do a one-to-many merge of the four datasets by year. The
266     grains dataset has up to four records for each year, so the other three
267     datasets will need to be replicated.;
268
269 * First, we must sort all data sets by year. This macro will sort an arbitrary
270     number of datasets. Note that it mutates currently existing datasets rather
271     than assigning new names to the sorted datasets.;
272
273 %MACRO SORTALL (DAT = , BYVAR = );
274     %LET N = %SYSFUNC(COUNTW(&DAT));
275     %DO I = 1 %TO &N;
276         PROC SORT DATA = %SCAN(&DAT, &I);
277             BY &BYVAR;
278         RUN;
279     %END;
280 %MEND;
281
282 %SORTALL(
283     DAT = ALLGRNS INFLATION PRES TEMP,
284     BYVAR = YEAR
285 );
286
287 * Now we can do the actual merge. Only the records with admissible years
288     (specified by the macro variables &MINYEAR and &MAXYEAR respectively)
289     will be read in and included in the merge.;
```

```
290
291  DATA HOME.GRAINS;
292      MERGE ALLGRNS INFLATION PRES TEMP;
293      WHERE &MINYEAR <= YEAR <= &MAXYEAR;
294      BY YEAR;
295  RUN;
296
297  PROC SORT DATA = HOME.GRAINS;
298      BY GRN YEAR;
299
```