

meat.sas: Explanation of code

Goals of code:

- Analyze a subset of the data
- Setup to calculate regression predictions at many X values
- Calculating confidence and prediction intervals
- Only printing a few observations
- Overlaying predicted values
- ANOVA lack of fit test

Analyzing a subset of the data: where or data set; if

SAS provides a couple of ways to restrict an analysis to a subset of the data. If you want to do many analyses or plots with a subset, I find it easier to create a second data set with the desired subset of data. If you want to do one analysis, it can be easier to specify that subset in the analysis.

The first proc glm illustrates specifying a subset in the analysis. This is done with a **where** statement. **where** is followed by a logical expression, e.g., **time <= 6**. The analysis will use only the rows of data where this is true. So here, only times 1, 2, 4, and 6. Not 8.

To create a new data set with a subset of values: This is done in a data step, **data meat16;**, followed by **set** to read observations from the meat data set. The crucial statement is **if** followed by a location expression. The new data set will only include observations where that logical statement is true.

Logical operators: SAS provides an extensive list of logical operators. These include

Symbol	meaning
=	equals
<	less than
>	greater than
<=	less than or equal
>=	greater than or equal
<>	not equal
in (1,3,6)	will be true if first argument matches any of the values in the ()

Setup to calculate regression predictions at many X values: data meat2;

The estimate statement demonstrated in meat1.sas prints predicted values for specified X values. Each new X requires a new estimate statement. If you want predicted values for more than a few observations, it is easier to add the desired X's to the data set.

The key idea is that if you provide an observation with an X value and a missing Y value, it won't be used to estimate parameters (because Y is missing) but you can calculate predicted values and related information (because X is specified). You just can't calculate a residual (again because Y is missing).

The two `data meat2;` steps show two ways to do similar things. You only need one of those data steps. The first data step (with the `datalines` statement) reads a set of new X values at which you want to make predictions. The code to indicate a missing numeric value in SAS is `.` (not in quotes). The string `'.'` in quotes is the character value `..`. The lonely `.` without quotes is the numeric missing value. This can be used in any data set. In particular, if a value for one variable is missing, use a lonely `.` to indicate that.

The second `data meat2;` step uses a loop to sequentially generate 1.0, 1.25, 1.5, ... 8. `ph` is set to a missing value. The `do` loop has the obvious syntax: `do variable = start to end by step.` If `by step` is omitted, the loop steps by 1, so `do i = 1 to 5` generates `i=1, 2, 3, 4, 5`. The `do` loop ends with the `end;` statement. All statements between `do` and `end;` are executed each time through the loop. Note: If you write `do time = 1, 5;`, SAS will execute the loop twice, once with `time = 1` and then with `time = 5`.

The loop sets the value of the time variable. In addition, we want to calculate the log transformed time, with is done the usual way. The `output;` writes the observation to the data set. An explicit `output;` statement is needed here because we are generating a sequence of observations, and we need to write each one to the data set.

The last piece of setup is to create a data set (to be called `all`) with both the real data (`meat`) and the prediction points (`meat2`). We have used the `set` command before to read observations from a pre-existing SAS data set. When there are multiple data set names, `set` will read all observations from the first data set, then all from the second, and then the third, and fourth, ... for all data sets. So, `set meat meat2;` will concatenate the `meat` and `meat2` data sets. We use that data set to fit the regression.

Calculating confidence and prediction intervals: `output;` statement

The `meat1.sas` program demonstrates fitting a regression using `proc glm`. The new thing here are the keywords in the `output` statement. The `output` statement can save many pieces of information. Each desired piece is specified by `keyword = name`. The keyword specifies what piece of information you want; the name is the variable name in which that piece is stored. You choose the variable name. I indicate some of the most useful keywords. Comments in the `meat.sas` file describe what each keyword gives you. A full list of keywords can be found in the SAS documentation for `proc glm`.

Only printing a few observations: `data = resid (obs=5)`

Sometimes, I want to print a few observations in a data set to get variable names (e.g., after a `proc import`) or check for gross errors. Anytime you explicitly indicate the name of a data set (by `data=`), you can add a data set option in `()`. `obs=5` says to only process the indicated number of observations, so `data=resid (obs=5)` prints the first 5 observations in the `resid` dataset. `data=resid (obs=25)`

prints the first 25 observations.

Note: You can use `obs=` in any proc step, but you probably don't want to. If you used it in a proc step that does a statistical analysis, that analysis would only use a subset of the data.

Overlaying observations and predicted values: `proc sgplot;`

We have previously used `proc sgplot; scatter ;` to plot data. You can provide multiple plot commands which will overlay all the requested pieces. The series command draws lines (by connecting the observations). The first `proc sgplot` overlays the data and fitted line. The second adds the lower and upper prediction limits. The `/lineattrs` option (line attributes) modifies how the line is drawn. The suboptions go in `()`, exactly as indicated here. Besides color and pattern, there is also a width suboption. There are many other options besides `lineattrs` and many other drawing commands besides `scatter` and `series`. The SAS documentation for `proc sgplot` gives you all the details.

ANOVA lack of fit test: `data meat3;` and subsequent `proc glm.`

To construct the ANOVA lack of fit test, we need to create a copy of the X variable. The `meat3` data step does that. The copy is called `ctime`, but the name is immaterial. It just has to be something other than the original name. The actual work is done by the `proc glm` with the `model ph = time ctime / ss1;` statement. One copy of the X variable goes in a class statement (to fit the ANOVA model with a separate mean for each unique X value). The model statement fits the regression (`time`) followed by the ANOVA model (`ctime`). The change in SS is the SS for lack of fit. The `/ss1` option requests only sequential SS (SAS type 1). That's because some of the type III SS don't make sense here (e.g. for `ctime` followed by `time`).

Note: If the first variable in the model is the one specified in the class statement), the SS and df for the second term are both 0. Fix by reversing the order.