This mini-document shows an example of a two-sample $t$-test done with SAS. The data are the "degree of reading power" data from Moore and McCabe.

First, the data. These are in a file `drp.dat` with each observation on one line. Some of the data file looks like this:

```
t 43
t 53
t 57
t 49
t 56
t 33
c 42
c 33
c 46
c 37
```

Each observation has two parts: first, an indication of whether it comes from the treatment or control group, and second, the value (score on a reading test).

A SAS program has two parts: first, a "data step" where you read in the data from the file. You have to tell SAS what the file is called, and then what the variable(s) are called. Also, if any of the variables are (or could be) characters rather than numbers, you have to tell SAS this as well. So our data step looks like this:

```
options linesize=80;

data drp;
  infile "drp.dat";
  input group $ score;
```

The first line just makes sure that the printed output doesn't go off the side of the page. Then comes the data step in earnest: `infile` to give the name of the data file (in quotes) and `input` to give the names of the variables, with `group` as a character variable, `t` or `c`. This is the easiest approach for reading in data. (SAS has a lot of variations on this, which are only worth learning if you need to know.)

This is followed by one or more "proc steps", which instruct SAS what to do with your data. There are many different "procs" that SAS knows about; usually each statistical procedure has its own proc. We'll use two procedures: first PROC MEANS to calculate the mean score for each group:

```
proc means;
  class group;
  var score;
```

then PROC TTEST to do a two-sample $t$-test to look for statistically significant differences between the two groups (ie. to answer the question "did the students in the treatment group become better readers than those in the control group?").

```
proc ttest;
  class group;
  var score;
```

Note, in each case, the PROC line is followed by a couple of other lines specifying how we want the PROC to run. In both cases here, we're saying that `group` is a classification variable (divides the data into groups), and `score` is what we're looking for differences in. Put the data step and the two proc steps together (one after the other) in a file called `drp.sas`, and then run SAS on a command line by typing `sas drp.sas`. There will be no output (apart maybe from a cryptic message), but there will be two new files, `drp.log` and `drp.lst`.

The log file `drp.log` tells you whether everything worked. If it did, you'll see an echo of the lines in `drp.sas` with some comments (about the number of lines in the data set, about the number of pages of output, and about how long the procedures took to run). If not, you'll get some kind of clue about why not by looking at the log file. Usually it's some kind of typo (one time, I typed `clsss group`, and SAS told me, via the log file, that it didn't understand `clsss`).

If everything ran successfully, you'll find the output in `drp.lst`. In this case, there are two "pages" to the output: one from PROC MEANS and one from PROC TTEST. PROC MEANS is straightforward enough:

```
                            The SAS System                                1
                                          23:19 Tuesday, September 30, 2008


                            The MEANS Procedure

                          Analysis Variable : score

           N
group     Obs    N          Mean        Std Dev        Minimum        Maximum
--------------------------------------------------------------------------------
c          23   23    41.5217391     17.1487332     10.0000000     85.0000000

t          21   21    51.4761905     11.0073568     24.0000000     71.0000000
--------------------------------------------------------------------------------
```

This gives you sample sizes, means, SDs, min and max for each group. You get a quick clue about whether the data were entered correctly, and eyeballing this suggests that the treatment students scored about 10 marks higher than the control students on average.

PROC TTEST produces a lot of stuff. This is the way SAS works: there is probably going to be more output than you need, and it's up to you to pick out what you want and ignore the rest. Here we go:

```
                            The SAS System                                2
                                          23:19 Tuesday, September 30, 2008
```

```
                          The TTEST Procedure

                              Statistics

                        Lower CL          Upper CL  Lower CL
     Variable  group      N    Mean    Mean    Mean   Std Dev  Std Dev

     score     c          23  34.106  41.522  48.937  13.263   17.149
     score     t          21  46.466  51.476  56.487  8.4213   11.007
     score     Diff (1-2)     -18.82  -9.954  -1.091  11.998   14.551

                              Statistics

                            Upper CL
         Variable  group     Std Dev   Std Err   Minimum   Maximum

         score     c          24.271    3.5758        10        85
         score     t          15.895     2.402        24        71
         score     Diff (1-2) 18.495    4.3919


                               T-Tests

      Variable    Method         Variances      DF   t Value   Pr > |t|

      score       Pooled         Equal          42    -2.27     0.0286
      score       Satterthwaite  Unequal      37.9    -2.31     0.0264


                        Equality of Variances

          Variable    Method    Num DF   Den DF   F Value   Pr > F

           score      Folded F      22       20      2.43   0.0507
```

The tables of Statistics give you sample means and stuff, like PROC MEANS
did (you can check that the sample means and SDs are the same). Also, there
are confidence intervals for the population means and SDs for each group singly
and for the two groups compared, which is what we want. The output tells us
that the 95% confidence interval for the difference between population means
is from –18.8 to –1.1 (control minus treatment). This says we are pretty sure
that the treatment group mean is higher, but the data don't tell us much about
how much higher. (If you want something other than a 95% interval, say 99%,
change your PROC TTEST line by adding `alpha=0.01` to the end of it.)

The tests in the next table give the two versions of the two-sample $t$-test.
The pooled test assumes that the two population SDs are equal. You assess
this by looking at the sample SDs, which are about 17 and 11: not equal, but
not so far apart. If you want, you can look at the bottom table, Equality of
Variances, which does a formal test. Since the P-value here of 0.0507 is pretty

small, it suggests that the assumption of equal population SDs is shaky. So using the Satterthwaite test is safer, though you can see that the P-values for the two tests, 0.0286 and 0.0264, are very close. Either way, the P-value for comparing the means is less than 0.05 (two-sided), and so we can conclude that the observed difference of 10 points is statistically significant, real, reproducible or whatever you will.