# 1   Topic

**PSPP, an alternative to SPSS.**

I spend a lot of time to analyze the available free statistical and data mining tools. There is not bad software, but some tools are more appropriate for some tasks. Thus, we must identify the one which is the best suited to our configuration. For that, we must know a large number of tools.

In this tutorial, we describe PSPP. It is presented as an alternative to the well-known SPSS: "*PSPP is a program for statistical analysis of sampled data. It is a free replacement for the proprietary program SPSS, and appears very similar to it with a few exceptions*[1]*"*. Instead of to describe in detail each feature, the documentation is available on the website, we present some statistical techniques. We compare the results with those of **Tanagra**, **R 2.13.2** and **OpenStat (build 24/02/2012)**. This is also a way to validate them. If they provide different results, it means that there is a problem.

# 2   Dataset

We use a version of the "Automobile" dataset from the UCI server[2]. According to the statistical method that we analyze, we use some variables of the dataset. The most important here is to show how to perform the various analyses with PSPP.

# 3   PSPP

## 3.1   Loading and installing PSPP



---

[1] http://www.gnu.org/software/pspp/pspp.html

[2] http://archive.ics.uci.edu/ml/datasets/Automobile

PSPP is downloadable online (http://www.gnu.org/software/pspp/). We use the **0.7.8** (2011/11/11) version. PSPP needs the MINGW[3] environment, but fortunately, the installation of this last one is done automatically. Thus, the installation process under Windows is performed easily.

## 3.2    Command line mode

We can set instructions in a program file (e.g. with a standard text editor) and send this one to the executable file PSPP.EXE. The syntax of the commands is the same to the one of SPSS. By learning to program with PSPP, we will know to do this with SPSS. This is really interesting.

For instance, we want to compare the horsepower of the cars according to their fuel type using a t-test. We set the following commands into the "**test.syn**" program file.

```
GET FILE="D:\dataset\pspp\autos.sav".


T-TEST /VARIABLES= horsepower
      /GROUPS=fuel_type("gas","diesel")   /MISSING=ANALYSIS
      /CRITERIA=CIN(0.95).
```

The results are displayed into the MSDOS console.



PSPP compares first the conditional variances with the Levene's test. Then, it compares the means with and without the homoscedasticity assumption.

---

[3] http://www.mingw.org/

### 3.3    Terminal mode

We can use also PSPP in a terminal mode. After we launch PSPP.EXE, a command interpreter allows us to set the instructions. The results are displayed into the same window. Here, we load the "autos.sav" data file, and we display the dataset dictionary.



### 3.4    Menu-driven mode



The easiest way to handle PSPP is the menu-guided mode. We use it in this tutorial. The features are grouped in menus "données" (dataset), "transformation", and "analyse" (analysis). *Curiously,*

*although that I have installed the English version, some menu-items are in French on my computer. I do not know if we have the same phenomenon with an English language operating system.*

# 4   Importing the data file

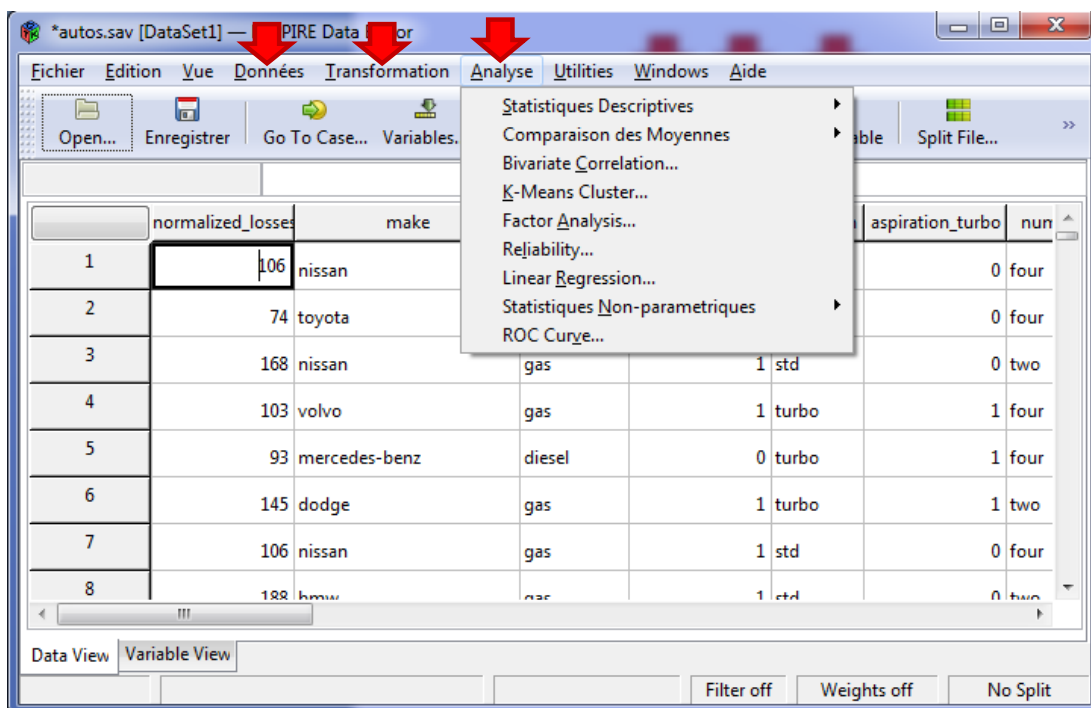First, we must import the « **autos_pspp.txt** » tab delimited text file. We click on the FICHIER / IMPORT DELIMITED TEXT DATA menu. We pick the file into the dialog box.

A wizard appears. We must: (1) import all the instances; (2) the first row corresponds to the variable name; (3) the TAB character is the column separator; (4) the values are string (chaîne) or numeric (numérique). We click on the APPLIQUER button.

We save the loaded dataset into the PSPP native format ("**autos.sav**"). This is the same one as SPSS.

# 5   A few statistical methods with PSPP

In this section, we present some statistical techniques available in PSPP. When this is possible, we compare the results with those of Tanagra.

## 5.1   Descriptive statistics – Numeric variables

To obtain the description for horsepower and city mpg, we click on the ANALYSE / STATISTIQUES DESCRIPTIVES / DESCRIPTIVES menu. Into the dialog settings, we select the indicators to compute.



The results are displayed into a new window ("Output viewer").

The MORE UNIVARIATE CONT STAT of **Tanagra** provides the same values.



| Attribute | Stats | |
|---|---|---|
| horsepower | **Statistics** | |
| | Average | 104.2537 |
| | Median | 95 |
| | Std dev. [Coef of variation] | 39.5192 [0.3791] |
| | MAD [MAD/STDDEV] | 30.2093 [0.7644] |
| | Min * Max [Full range] | 48.00 * 288.00 [240.00] |
| | 1st * 3rd quartile [Range] | 70.00 * 116.00 [46.00] |
| | Skewness (std-dev) | 1.3980 (0.1698) |
| | Kurtosis (std-dev) | 2.6785 (0.3381) |
| | | |
| city_mpg | **Statistics** | |
| | Average | 25.2195 |
| | Median | 24 |
| | Std dev. [Coef of variation] | 6.5421 [0.2594] |
| | MAD [MAD/STDDEV] | 5.2155 [0.7972] |
| | Min * Max [Full range] | 13.00 * 49.00 [36.00] |
| | 1st * 3rd quartile [Range] | 19.00 * 30.00 [11.00] |
| | Skewness (std-dev) | 0.6637 (0.1698) |
| | Kurtosis (std-dev) | 0.5786 (0.3381) |
| | | |

## 5.2   Conditional descriptive statistics

We want calculate the descriptive statistics for horsepower according to the values of fuel type ("gas" or "diesel").

We click on the ANALYSE / STATISTIQUES DESCRIPTIVES / EXPLORE menu. We set "horsepower" into DEPENDENT LIST, "fuel type" into FACTOR LIST. By clicking the STATISTICS button, we can specify the indicators to compute. We validate.

For instance, we observe that the mean of horsepower is 84.45 for the cars corresponding to "fuel type = gas", and 106.39 for "fuel type = diesel". We obtain very detailed results.

The GROUP CHARACTERIZATION component of **Tanagra** provides the same values, but the results are definitely less detailed.



## 5.3   Contingency table

We want to analyze the relation between fuel type and aspiration ("turbo" or "standard") with a cross tab. We click on the ANALYSE / STATISTIQUES DESCRIPTIVES / CROSSTABS menu. We set the first variable into ROWS list, the second one into the COLUMNS.

fuel_type * aspiration [Compter, ligne %, colonne %, total %].

|  | aspiration | | |
|---|---|---|---|
| fuel_type | std | turbo | Total |
| diesel | 7.0 | 13.0 | 20.0 |
|  | 35.0% | 55.0% | 100.0% |
|  | 4.2% | 35.1% | 9.8% |
|  | 3.4% | 5.3% | 9.8% |
| gas | 161.0 | 24.0 | 185.0 |
|  | 37.0% | 13.0% | 100.0% |
|  | 95.8% | 54.9% | 90.2% |
|  | 78.5% | 11.7% | 90.2% |
| Total | 168.0 | 37.0 | 205.0 |
|  | 82.0% | 18.0% | 100.0% |
|  | 100.0% | 100.0% | 100.0% |
|  | 82.0% | 18.0% | 100.0% |

Tests du Chi-Deux

| Statistique | Valeur | df | Asymp. Sig. (2-tailed) |
|---|---|---|---|
| Pearson Chi-Square | 33.03 | 1 | .00 |
| Likelihood Ratio | 24.90 | 1 | .00 |
| Continuity Correction | 29.61 | 1 | .00 |
| N observations valides | 205 | | |

Here also, PSPP can provide very detailed results. Some measures of association (Theil's U, Cohen's Kappa, etc.) and the various percentages are also displayed.

The same results are available under **Tanagra** with the CONTINGENCY CHI-SQUARE component. But the organization is a little different.



| Row (Y) | Column (X) | Statistical indicator | | Cross-tab | | | |
|---|---|---|---|---|---|---|---|
|  |  | Stat | Value | | std | turbo | Sum |
|  |  | d.f. | 1 | gas | 161 | 24 | 185 |
|  |  | Tschuprow's t | 0.401397 | | 87.03% | 12.97% | 100% |
|  |  | Cramer's v | 0.401397 | diesel | 7 | 13 | 20 |
| fuel_type | aspiration | Phi² | 0.16112 | | 35.00% | 65.00% | 100% |
|  |  | Chi² (p-value) | 33.03 (0.0000) | Sum | 168 | 37 | 205 |
|  |  | Lambda | 0 | | 82% | 18% | 100% |
|  |  | Tau (p-value) | 0.1611 (0.0000) | | | | |
|  |  | U(R/C) (p-value) | 0.1900 (0.0000) | | | | |

TARGET : (fuel_type)
INPUT : (aspiration)

## 5.4  Comparison of two means – Independent samples

Beyond the conditional descriptive statistics, we can test if the means are significantly different. We click on the ANALYSE / COMPARAISON DES MOYENNES / INDEPENDENT SAMPLES T TEST, we set horsepower into TEST VARIABLE list, fuel type into DEFINE GROUPS.

PSPP compares the conditional variances first by using the Levene's test. We observe that they are not significantly different at the 5% level (p-value = 0.17). However, it displays the comparison of means with and without the homoscedasticity assumption. In both cases, we see that the means are significantly different.

We use three components under **Tanagra**. But they are plugged after the same DEFINE STATUS, which specify the role of the variables, into the diagram. There is not repetitive handling.



Of course, the results are identical to those of PSPP.

## 5.5    Comparison of two means – Paired samples

Now, we want to compare the city mpg and the highway mpg. The difference has to be computed for each car i.e. we have paired samples. We click on the ANALYSE / COMPARAISON DES MOYENNES / PAIRED SAMPLES T TEST menu. We select the pair of variables "city mpg" and "highway mpg".

PSPP provides the mean for each variable (25.2 city mpg, 30.75 highway mpg), the correlation between them, and the details for the statistical test. Patently, the consumption is higher in city (the miles that we can cover with one gallon of fuel is significantly lower).



Le PAIRED T-TEST component of **Tanagra** provides the same values.



## 5.6    Comparison of K means – Analysis of variance (ANOVA)

We want now to compare the means of several (> 2) groups. We click on the ANALYSE / COMPARAISON DES MOYENNES / ONE WAY ANOVA menu.

We set "horsepower" as DEPENDENT VARIABLE, and "body style" as FACTOR. PSPP assess the homoscedasticity assumption with the Levene's test. Then, it performs the comparison of means. At the 5% level, we reject the null hypothesis.

**Tanagra** provides the same results. But we use two components.



## 5.7    Multiple regression

We want to explain the consumption (city mpg) from the fuel type (dummy variable), the aspiration, the curb weight and the horsepower. We set these parameters into the dialog box of ANALYSE / LINEAR REGRESSION menu.



PSPP provides the multiple correlation coefficient R; the coefficient of determination R-squared; the ANOVA table for the regression, the table of coefficients. The model is globally significant at the 5% level. And all the coefficients are significant. "Weight" and "fuel type" have the highest influence on the consumption.

The MULTIPLE LINEAR REGRESSION of **Tanagra** provides the same results.

## Global results

| | |
|---|---|
| Endogenous attribute | city_mpg |
| Examples | 205 |
| R² | 0.790990 |
| Adjusted-R² | 0.786810 |
| Sigma error | 3.020670 |
| F-Test (4,200) | 189.2233 (0.000000) |

TARGET : (city_mpg)
INPUT : (fuel_type_gas,
aspiration_turbo,
curb_weight, horsepower)

└─ Define status 7
   └─ Multiple linear regression 1

## Analysis of variance

| Source | xSS | d.f. | xMS | F | p-value |
|---|---|---|---|---|---|
| Regression | 6906.2327 | 4 | 1726.5582 | 189.2233 | 0.0000 |
| Residual | 1824.8892 | 200 | 9.1244 | | |
| Total | 8731.1220 | 204 | | | |

## Coefficients

| Attribute | Coef. | std | t(200) | p-value |
|---|---|---|---|---|
| Intercept | 57.640300 | 1.681678 | 34.275470 | 0.000000 |
| fuel_type_gas | -8.597166 | 0.929372 | -9.250509 | 0.000000 |
| aspiration_turbo | -1.641124 | 0.636844 | -2.576963 | 0.010687 |
| curb_weight | -0.007887 | 0.000722 | -10.917917 | 0.000000 |
| horsepower | -0.040383 | 0.009654 | -4.183196 | 0.000043 |

## 5.8   ROC curve

Cars are initially assigned a risk factor symbol associated with its price. Then, if it is more risky (or less), this symbol is adjusted by moving it up (or down) the scale. Insurance companies call this process "symboling"[4]. From this variable, we define a new column "risky". Its value is "yes" if "symboling > 0", "no" otherwise. At the same time, the "normalized losses" is defined as the relative average loss payment per insured vehicle year. This value is normalized for all autos within a particular size classification (two-door small, station wagons, sports/specialty, etc...).

We were wondering if the "normalized losses" allows to distinguish risky cars from non-risky ones. We want to use the ROC curve for answering to this question. We click on the ANALYSE / ROC CURVE menu. We set the following settings.



The "positive" instances corresponds to "risky = yes" values. PSPP calculates automatically the AUC criterion (area under curve).

---

[4] http://archive.ics.uci.edu/ml/datasets/Automobile

Case Summary

| risky | Valid N (listwise) | |
|---|---|---|
| | Unweighted | Weighted |
| Positive | 113 | 113.00 |
| Negative | 92 | 92.00 |

Area Under the Curve (normalized_losses)

| Area |
|---|
| .74 |

There are 113 positive instances into the dataset. AUC = 74% i.e. the probability that a randomly chosen risky cars has higher normalized losses than a randomly chosen non-risky cars is 74%.

For this time, **Tanagra** seems not agree. The shape of the ROC curve is a little different. In effect, Tanagra cuts the score value in 20 intervals and build the ROC curve from the corresponding values. Thus, we obtain **a smooth curve**. But the underlying values used for the construction of the curve are the same. We obtain the same AUC criterion.



TARGET : (city_mpg)
INPUT : (fuel_type_gas, aspiration_turbo, curb_weight, horsepower)

Sample size : 205
Positive examples : 113
Negative examples : 92

| Score Attribute | normalized_losses | | |
|---|---|---|---|
| AUC | 0.7461 | | |
| Target size (%) | Score | FP-Rate | TP-Rate |
| 0 | 256.0000 | 0.0000 | 0.0000 |
| 5 | 188.0000 | 0.0326 | 0.0619 |
| 10 | 161.0000 | 0.0543 | 0.1327 |
| 15 | 158.0000 | 0.1087 | 0.1770 |
| 20 | 150.0000 | 0.1087 | 0.2743 |
| 25 | 137.0000 | 0.1196 | 0.3540 |
| 30 | 129.0000 | 0.1304 | 0.4336 |
| 35 | 122.0000 | 0.1413 | 0.5133 |
| 40 | 122.0000 | 0.1848 | 0.5752 |
| 45 | 122.0000 | 0.2391 | 0.6195 |
| 50 | 122.0000 | 0.2717 | 0.6814 |
| 55 | 122.0000 | 0.3370 | 0.7168 |
| 60 | 115.0000 | 0.3913 | 0.7699 |
| 65 | 106.0000 | 0.4674 | 0.7965 |
| 70 | 103.0000 | 0.5000 | 0.8584 |
| 75 | 101.0000 | 0.5870 | 0.8761 |
| 80 | 94.0000 | 0.6413 | 0.9292 |
| 85 | 91.0000 | 0.7283 | 0.9469 |
| 90 | 85.0000 | 0.8152 | 0.9646 |
| 95 | 78.0000 | 0.8913 | 0.9912 |
| 100 | 65.0000 | 1.0000 | 1.0000 |

# 6   Perform the same analyses with other tools

Of course, we can implement the same statistical approaches with other tools. In this section, we show briefly the whole diagram under **Tanagra**, we details the commands under **R 2.13.2**, and we show the available features of **OpenStat** (which is really similar to PSPP).

## 6.1   Tanagra

Some statistical approaches available in PSPP are not available in Tanagra, and conversely. For the methods described in this tutorial, here is the whole diagram under Tanagra.



## 6.2   R software

Excluding the ROC curve, we detail here the commands and the corresponding outputs in R. In some circumstances, we need a specific package that we load with the **library(.)** instruction.

```
> #loading the dataset
> setwd("D:/DataMining/Databases_for_mining/logiciels_dataset/pspp")
> autos <- read.table(file="autos_pspp.txt",header=T,sep="\t",dec=".")
>
> #descriptive statistics
> print(summary(data.frame(autos$horsepower,autos$city_mpg)))
 autos.horsepower autos.city_mpg
 Min.   : 48.0    Min.   :13.00
 1st Qu.: 70.0    1st Qu.:19.00
 Median : 95.0    Median :24.00
 Mean   :104.3    Mean   :25.22
 3rd Qu.:116.0    3rd Qu.:30.00
```

```
 Max.   :288.0    Max.    :49.00
>
> #conditionnal descriptive statistics
>
print(tapply(autos$horsepower,autos$fuel_type,FUN=function(x){c(m=mean(x),s
=sd(x))}))
$diesel
       m         s
84.45000 25.95842

$gas
        m         s
106.39459  40.18342

> #crosstabs and test of independance
> library(gmodels)
>
print(CrossTable(autos$fuel_type,autos$aspiration,prop.r=F,prop.c=F,prop.t=
F,chisq=T))

   Cell Contents
|-----------------------|
|                     N |
| Chi-square contribution |
|-----------------------|

Total Observations in Table:  205


             | autos$aspiration
autos$fuel_type |       std |     turbo | Row Total |
---------------|-----------|-----------|-----------|
       diesel |         7 |        13 |        20 |
             |     5.380 |    24.427 |           |
---------------|-----------|-----------|-----------|
          gas |       161 |        24 |       185 |
             |     0.582 |     2.641 |           |
---------------|-----------|-----------|-----------|
   Column Total |       168 |        37 |       205 |
---------------|-----------|-----------|-----------|


Statistics for All Table Factors

Pearson's Chi-squared test
------------------------------------------------------------
Chi^2 =  33.02955    d.f. = 1     p =  9.076896e-09

Pearson's Chi-squared test with Yates' continuity correction
------------------------------------------------------------
Chi^2 =  29.60576    d.f. = 1     p =  5.294738e-08

$t
       y
x        std turbo
  diesel  7    13
  gas    161    24

$prop.row
       y
x            std      turbo
  diesel 0.3500000 0.6500000
  gas    0.8702703 0.1297297
```

```
$prop.col
        y
x                std      turbo
  diesel 0.04166667 0.35135135
  gas    0.95833333 0.64864865

$prop.tbl
        y
x                std      turbo
  diesel 0.03414634 0.06341463
  gas    0.78536585 0.11707317

$chisq

        Pearson's Chi-squared test

data:  t
X-squared = 33.0295, df = 1, p-value = 9.077e-09

$chisq.corr

        Pearson's Chi-squared test with Yates' continuity correction

data:  t
X-squared = 29.6058, df = 1, p-value = 5.295e-08

> #Levene test for variance homogeneity
> library(lawstat)
> print(levene.test(autos$horsepower,autos$fuel_type,location="mean"))

        classical Levene's test based on the absolute deviations from the
mean ( none not applied
        because the location is not set to median )

data:  autos$horsepower
Test Statistic = 1.9242, p-value = 0.1669

> #t-test for independent samples
> print(t.test(autos$horsepower ~ autos$fuel_type, var.equal=T))

        Two Sample t-test

data:  autos$horsepower by autos$fuel_type
t = -2.3861, df = 203, p-value = 0.01795
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -40.078454  -3.810736
sample estimates:
mean in group diesel    mean in group gas
          84.4500              106.3946

> #Welch t-test for independent samples
> print(t.test(autos$horsepower ~ autos$fuel_type, var.equal=F))

        Welch Two Sample t-test

data:  autos$horsepower by autos$fuel_type
t = -3.3693, df = 29.912, p-value = 0.00209
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
```

```
 -35.247706  -8.641483
sample estimates:
mean in group diesel    mean in group gas
            84.4500             106.3946

> #t-test for paired samples
> print(t.test(autos$city_mpg,autos$highway_mpg, paired=T))

        Paired t-test

data:  autos$city_mpg and autos$highway_mpg
t = -48.1901, df = 204, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -5.758033 -5.305382
sample estimates:
mean of the differences
           -5.531707

> #Levene test for variance homogeneity
> print(levene.test(autos$horsepower,autos$body_style,location="mean"))

        classical Levene's test based on the absolute deviations from the
mean ( none not applied
        because the location is not set to median )

data:  autos$horsepower
Test Statistic = 1.6904, p-value = 0.1536

> #analysis of variance
> print(aov(horsepower ~ body_style, data = autos))
Call:
   aov(formula = horsepower ~ body_style, data = autos)

Terms:
                body_style Residuals
Sum of Squares    17744.68 300856.13
Deg. of Freedom          4        200

Residual standard error: 38.78506
Estimated effects may be unbalanced
>
> #linear regression
>                    print(summary(lm(city_mpg                        ~
fuel_type_gas+aspiration_turbo+curb_weight+horsepower, data=autos)))

Call:
lm(formula = city_mpg ~ fuel_type_gas + aspiration_turbo + curb_weight +
    horsepower, data = autos)

Residuals:
    Min      1Q  Median      3Q     Max
-9.1931 -1.4955 -0.1292  0.8772 15.8097

Coefficients:
                  Estimate Std. Error t value Pr(>|t|)
(Intercept)      57.6402999  1.6816779  34.275  < 2e-16 ***
fuel_type_gas    -8.5971662  0.9293722  -9.251  < 2e-16 ***
aspiration_turbo -1.6411239  0.6368442  -2.577   0.0107 *
curb_weight      -0.0078871  0.0007224 -10.918  < 2e-16 ***
```

```
horsepower         -0.0403830  0.0096536  -4.183  4.3e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.021 on 200 degrees of freedom
Multiple R-squared: 0.791,      Adjusted R-squared: 0.7868
F-statistic: 189.2 on 4 and 200 DF,  p-value: < 2.2e-16
```
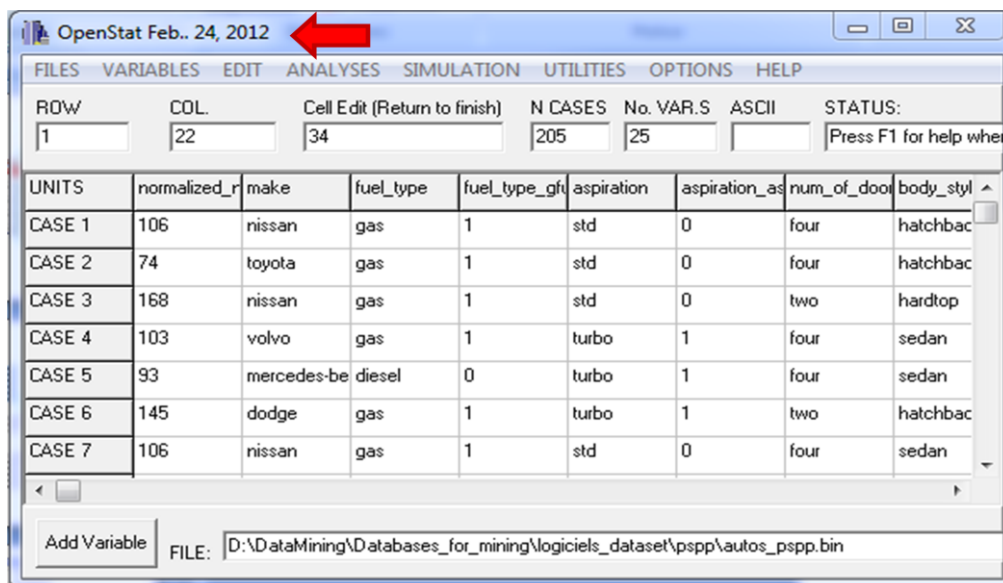
Knowing the appropriate commands for each method is the main difficulty under R. Fortunately, some websites provide a valuable assistance (e.g. Quick-R).

## 6.3  OpenStat

**OpenStat** is also a credible alternative to SPSS. In a previous tutorial, we have studied one of its variations (LazStat) in the regression analysis framework[5].
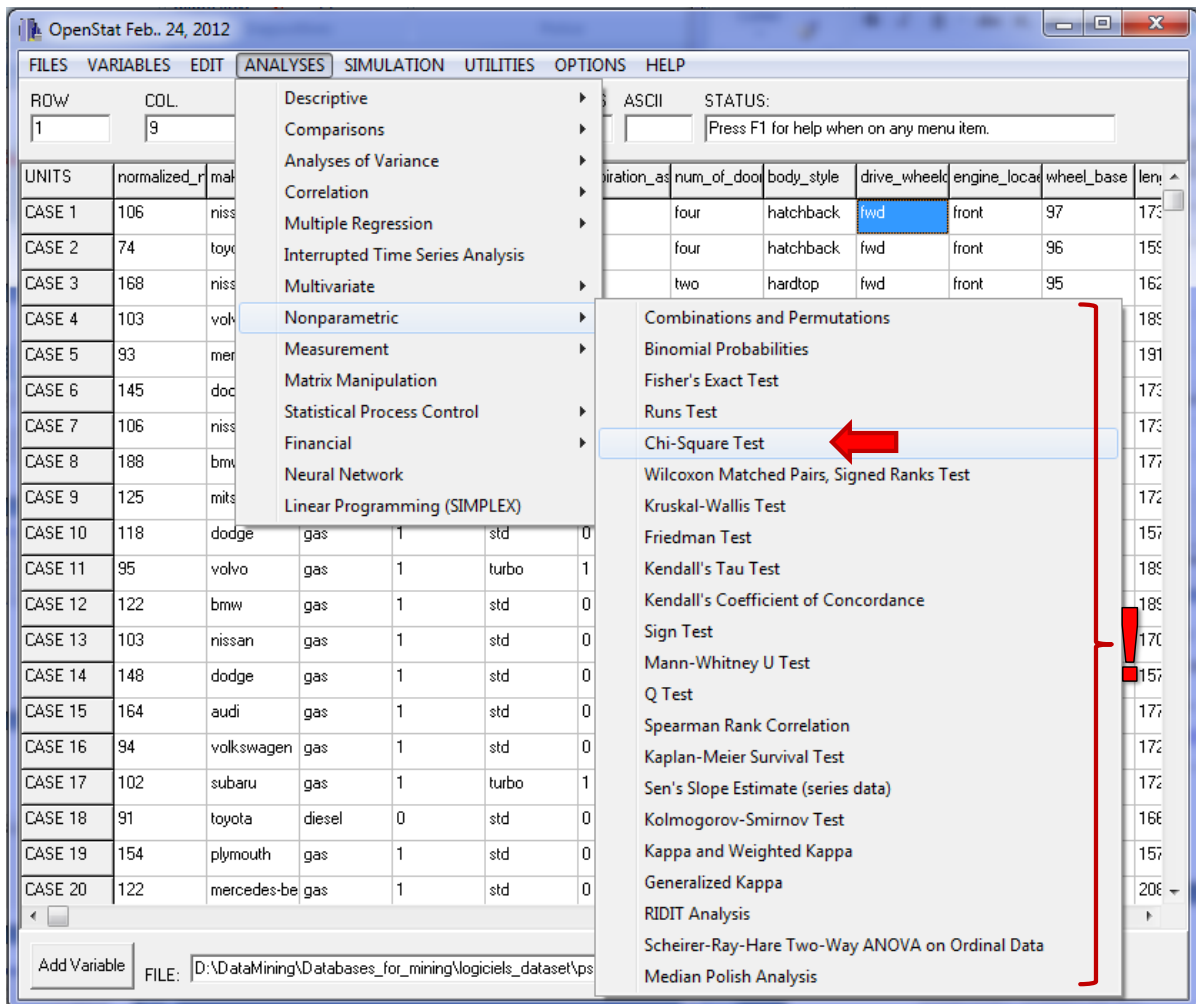
Like PSPP, we must import the data file first (FILES / IMPORT TAB FILE menu). The dataset is displayed into the grid.



We do not detail the results for each of the methods studied above. I perform these analyses elsewhere, I note that the results are the same as PSPP (as Tanagra and R). We describe only the results for the contingency table. We note that we must code the categorical variables as dummy ones before performing the analysis.

The statistical methods are available into the ANALYSES menu. We select the NONPARAMETRIC item. The list of available techniques is really long. We will try to describe them in a forthcoming tutorial.

---

[5] http://data-mining-tutorials.blogspot.com/2012/03/regression-analysis-with-lazstats.html

For the CHI-SQUARE TEST, we set the following settings and we click on COMPUTE.



A results window appears. **OpenStat** provides additional statistics such as the Mantel-Haenszel test of linear association, the coefficient of contingency, etc.

```
Results Window                                                                                    x

No. of Cases = 205

OBSERVED FREQUENCIES

              Frequencies
                COL. 1      COL. 2      Total

      Row 1          7          13          20
      Row 2        161          24         185
      Total        168          37         205

CHI-SQUARED VALUE FOR CELLS

              Chi-square Values
                COL. 1      COL. 2
Row 1            5.380      24.427
Row 2            0.582       2.641

Chi-square =    33.030 with D.F. = 1. Prob. > value =     0.000

Liklihood Ratio =    24.904 with prob. > value = 0.0000

G statistic =    24.904 with prob. > value = 0.0000

phi correlation = 0.4014

Pearson Correlation r = -0.4014

Mantel-Haenszel Test of Linear Association =    32.868 with probability > value = 0.0000

The coefficient of contingency =    0.373

Cramer's V =    0.401
```

Unlike to PSPP, because **OpenStat** is only a menu-guided program, it is not possible to store a description of the treatments. So, it is not easy to perform the same sequence of analyses if we have a new version of the dataset for instance.

# 7  Conclusion

PSPP is a promising project. The structure of the software is really well thought out. Each menu action is translated into a PSPP instruction. Thus, we can save the sequence of commands into a script file. For instance, for the analysis of variance described above (section 5.6), PSPP generates the following instruction.

```
ONEWAY /VARIABLES= horsepower BY body_style
        /STATISTICS=HOMOGENEITY .
```

PSPP already proposes a large part of the statistical methods. It will be further complemented in the future. This is a tool that I will follow with interest.