

《统计计算与SAS软件》

实验9 描述性统计量的计算、图示

1907402030 熊雄

2021年12月9日

1 实验目的

数据的简单加工并掌握`means`、`univariate`、`freq`过程。

2 实验内容

以`sashelp.class`数据集为例。

1. 用`means`过程计算`weight`的均值，标准差，极差，四分位极差，方差，偏度，峰度，并将这些统计量保存到数据集`resultmeans`中；
2. 用`univariate`过程计算`height`的均值，标准差，极差，中位数，众数，方差，偏度，峰度，并将这些统计量的值保存到数据集`resultuni`中；
3. 对变量`height`画正态概率图（Q-Q图）、茎叶图、盒形图；
4. 用`freq`过程统计男，女生的频数，并考虑性别与年龄是否独立？
5. 将性别变量的值加标签：`M`表示成“男”，`F`表示成“女”；
6. 给变量加中文标签：`name`的标签为“姓名”，`sex`的标签为“性别”，`age`的标签为“年龄”，`height`的标签为“身高”，`weight`的标签为“体重”；
7. 将`weight`分成三组，`<85`为一组，用“轻”表示，`85-100`为一组，用“中”表示，`>100`为一组，表示为“重”。并考虑`sex`与这种分组间是否独立；
8. 画`height`和`weight`两变量的散点图，数据点用“红色三角形”表示；
9. 画`height`的直方图、年龄的饼图（可以设置值，填充颜色、类型等）；
10. 按性别分组，计算不同年龄的平均身高。

3 代码实现

3.1 利用 *MEAN* 过程分析 *weight*

以下代码实现了用 *means* 过程计算 *weight* 的均值(mean), 标准差(std), 极差(range), 四分位极差(qrange), 方差(variance/var), 偏度(skewness/skew), 峰度(kurtosis/kurt), 并将这些统计量保存到数据集 *resultmeans* 中。

在 SAS 中提交如下代码:

```

1  PROC MEANS  data = sashelp.class mean std range qrange var skew
    kurt;
2  var weight;
3  label mean = '均值'
4        std = '标准差'
5        range = '极差'
6        qrange = '四分位数极差'
7        var = '方差'
8        skew = '偏度'
9        kurt = '峰度';
10 output out = work.resultmeans
11       mean = mean
12       std = std
13       range = range
14       qrange = qrange
15       var = var
16       skew = skew
17       kurt = kurt; /*将输出结果保存在数据集resultmeans*/
18 RUN;
```

可以得到如下输出结果:

SAS 系统						
MEANS PROCEDURE						
分析变量: Weight 体重 (磅)						
均值	标准差	极差	四分位数极差	方差	偏度	峰度
100.0263158	22.7739335	99.5000000	28.5000000	518.6520468	0.1833510	0.6833648

3.2 利用 *UNIVARIATE*过程分析*weight*

以下代码实现了用*univariate*过程计算*height*的均值，标准差，极差，中位数，众数，方差，偏度，峰度，并将这些统计量的值保存到数据集*resultuni*中。

在 *SAS* 中提交如下代码：

```
1 PROC UNIVARIATE data = sashelp.class;
2     var height;
3     output out = work.resultuni;
4 run;
```

可以得到如下输出结果：

SAS 系统			
UNIVARIATE PROCEDURE			
变量: Height (身高 (英寸))			
矩			
N	19	权重总和	19
均值	62.3368421	观测总和	1184.4
标准差	5.12707525	方差	26.2869006
偏度	-0.2596696	峰度	-0.1389692
未校平方和	74304.92	校正平方和	473.164211
变异系数	8.22479143	标准误差均值	1.17623173

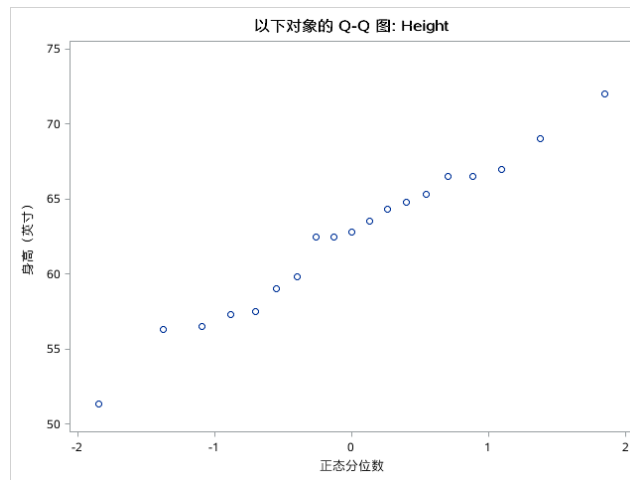
3.3 对变量*height*画图

3.3.1 对变量*height*画正态概率图（Q-Q图）

在 *SAS* 中提交如下代码：

```
1 PROC UNIVARIATE data = sashelp.class;
2     qqplot height;
3 run;
```

可以得到如下Q-Q图：



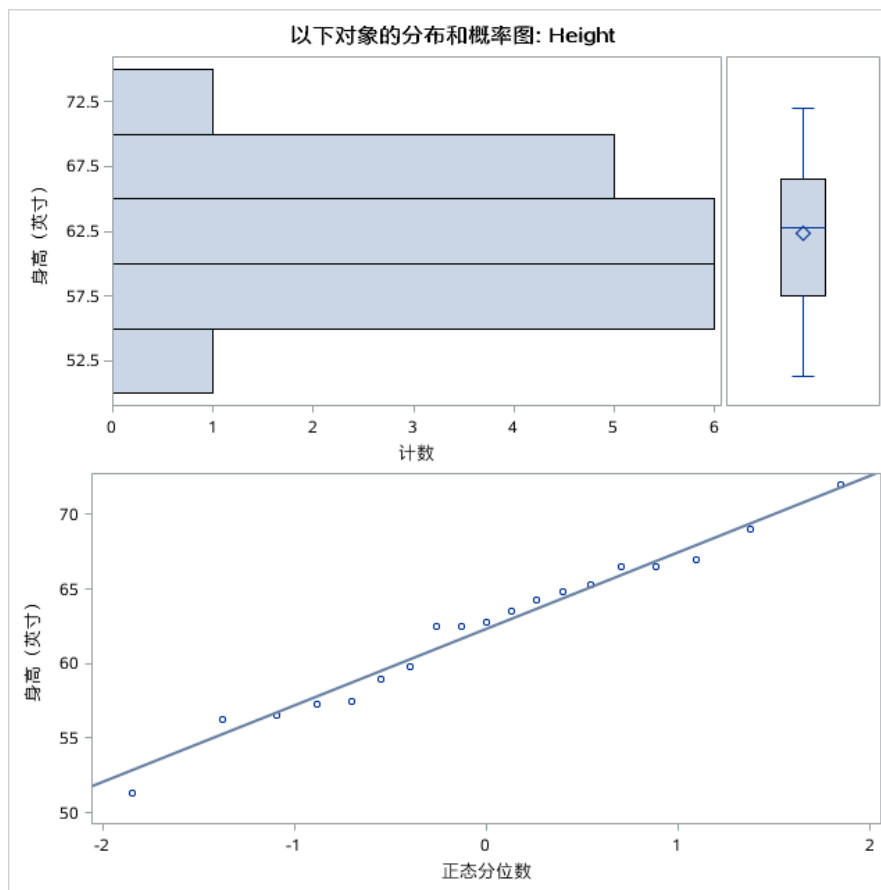
3.3.2 对变量 *height* 画茎叶图与盒形图

使用 PLOT 选项 UNIVARIATE 过程将产生三种图形：平行条状图 (Horizontal Bar Chart)、盒状图 (Box Plot)、正态分布拟合图 (Normal Probability Plot)。在 SAS 中提交如下代码：

```

1 ods select plots;
2 PROC UNIVARIATE data = sashelp.class plot;
3     var height;
4 RUN;
```

可以得到如下图形：



3.4 利用 *freq* 过程统计男, 女生的频数, 并考虑性别与年龄是否独立

在 *SAS* 中提交如下代码:

```
1 PROC FREQ data = sashelp.class;
2     table sex;
3     table sex * age /chisq;
4 RUN;
```

可以得到如下结果:

SAS 系统				
FREQ 过程				
性别				
Sex	频数	百分比	累积频数	累积百分比
男	10	52.63	10	52.63
女	9	47.37	19	100.00

频数 百分比 行百分比 列百分比	表 - Sex * Age							
	Sex(性别)	Age(年龄)						合计
		11	12	13	14	15	16	
		1	3	1	2	2	1	
男		5.26	15.79	5.26	10.53	10.53	5.26	52.63
		10.00	30.00	10.00	20.00	20.00	10.00	
		50.00	60.00	33.33	50.00	50.00	100.00	
女		5.26	10.53	10.53	10.53	10.53	0.00	47.37
		11.11	22.22	22.22	22.22	22.22	0.00	
		50.00	40.00	66.67	50.00	50.00	0.00	
合计		2	5	3	4	4	1	19
		10.53	26.32	15.79	21.05	21.05	5.26	100.00

表 "Age-Sex" 的统计量			
统计量	自由度	值	概率
卡方	5	1.4848	0.9148
似然比卡方检验	5	1.8748	0.8662
Mantel-Haenszel 卡方	1	0.0672	0.7955
Phi 系数		0.2795	
列联系数		0.2692	
Cramer V		0.2795	
WARNING: 100% 的单元格的期望计数 小于 5。卡方可能不是有效的检验。			
样本大小 = 19			

因此性别与年龄并不是相互独立的。

3.5 将性别变量的值加标签

利用 *FORMAT* 过程将性别变量的值加标签：*M* 表示成“男”，*F* 表示成“女”。在 *SAS* 中提交如下代码即可：

```

1 PROC FORMAT;
2     value $setsex 'M' = '男' 'F' = '女';
3 RUN;
4 DATA sss;
5     set sasHELP.class;
6     format sex setsex.;
7 run;

```

打开数据集 *work.sss*, 可以看到性别已经变成中文了:

	Name	Sex	Age	Height	Weight
1	Alfred	男	14	69	112.5
2	Alice	女	13	56.5	84
3	Barbara	女	13	65.3	98
4	Carol	女	14	62.8	102.5
5	Henry	男	14	63.5	102.5
6	James	男	12	57.3	83
7	Jane	女	12	59.8	84.5
8	Janet	女	15	62.5	112.5
9	Jeffrey	男	13	62.5	84
10	John	男	12	59	99.5
11	Joyce	女	11	51.3	50.5
12	Judy	女	14	64.3	90
13	Louise	女	12	56.3	77
14	Mary	女	15	66.5	112
15	Philip	男	16	72	150
16	Robert	男	12	64.8	128
17	Ronald	男	15	67	133
18	Thomas	男	11	57.5	85
19	William	男	15	66.5	112

3.6 给变量加中文标签

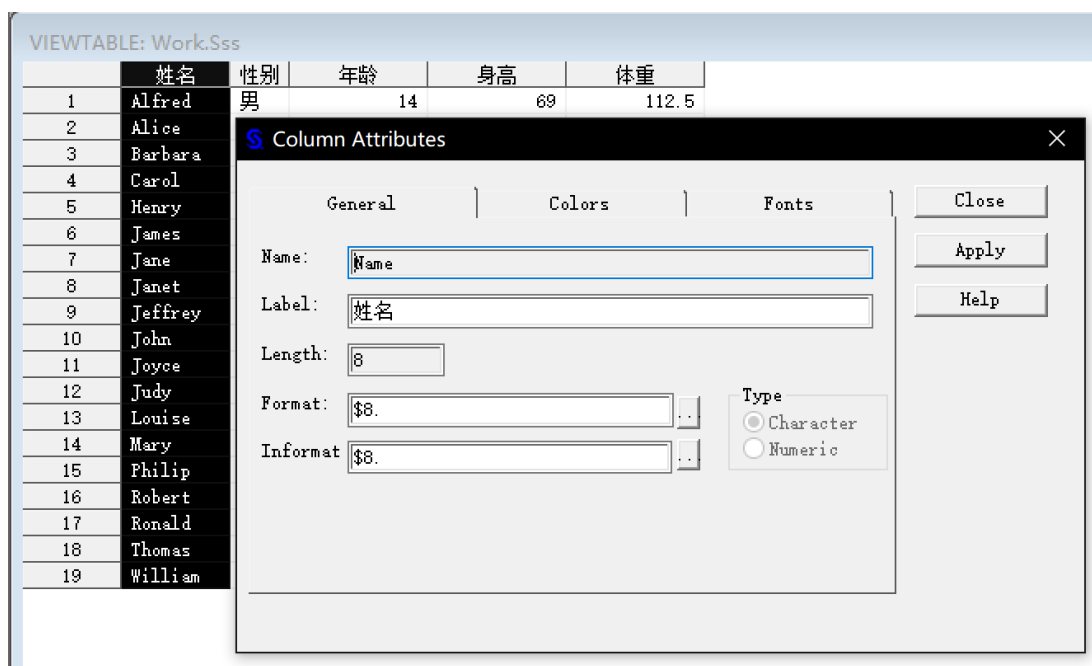
在 *DATA* 步中给变量加中文标签: *name* 的标签为“姓名”, *sex* 的标签为“性别”, *age* 的标签为“年龄”, *height* 的标签为“身高”, *weight* 的标签为“体重”。在 *SAS* 中提交如下代码:

```

1 DATA sss;
2 set sss;
3 label name = '姓名'
4     sex = '性别'
5     age = '年龄'
6     height = '身高'
7     weight = '体重';
8 run;

```

打开数据集`work.sss`，点开变量就可以看到Label了：



3.7 将`weight`分组

先利用`FORMAT`过程将`weight`分成三组的格式，<85为一组，用“轻”表示，85-100为一组，用“中”表示，>100为一组，表示为“重”，再利用`FREQ`过程实现分组。

```

1  PROC FORMAT;
2  VALUE setweight LOW - 85 = '轻'
3      85 - 100 = '中'
4      100 - HIGH = '重';
5  RUN;
6
7  PROC FREQ data = work.sss;
8      table weight;
9      format weight setweight.;
10     table sex * weight / chisq;
11 RUN;
```

得到输出如下：

SAS 系统

FREQ 过程

体重				
Weight	频数	百分比	累积 频数	累积 百分比
轻	7	36.84	7	36.84
中	3	15.79	10	52.63
重	9	47.37	19	100.00

频数
百分比
行百分比
列百分比

表 - Sex * Weight				
Sex(性别)	Weight(体重)			
	轻	中	重	合计
F	4	2	3	9
	21.05	10.53	15.79	47.37
	44.44	22.22	33.33	
	57.14	66.67	33.33	
M	3	1	6	10
	15.79	5.26	31.58	52.63
	30.00	10.00	60.00	
	42.86	33.33	66.67	
合计	7	3	9	19
	36.84	15.79	47.37	100.00

表 “Weight-Sex” 的统计量

统计量	自由度	值	概率
卡方	2	1.4275	0.4898
似然比卡方检验	2	1.4499	0.4844
Mantel-Haenszel 卡方	1	0.6384	0.4243
Phi 系数		0.2741	
列联系数		0.2644	
Cramer V		0.2741	
WARNING: 100% 的单元格的期望计数 小于 5。卡方可能不是有效的检验。			

样本大小 = 19

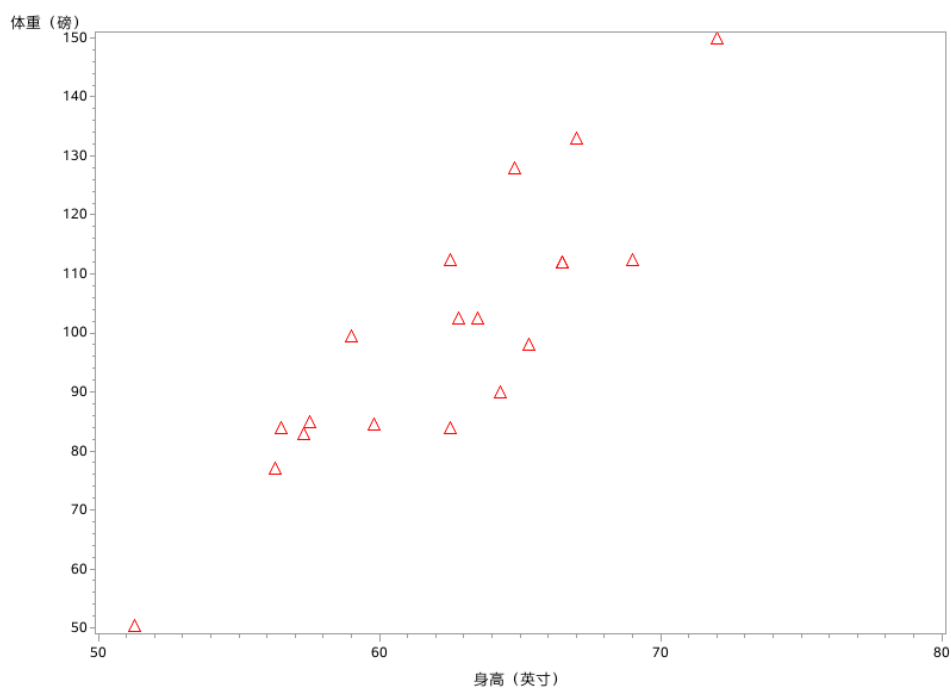
由“表‘Weight-Sex’的统计量”可以看到sex与这种分组间并不独立。

3.8 画 $height$ 和 $weight$ 两变量的散点图

画 $height$ 和 $weight$ 两变量的散点图，数据点用“红色三角形”表示；

```
1 proc gplot data = sashelp.class;
2     symbol V = TRIANGLE, CV = RED, H = 2;
3     plot weight * height;
4 run;
```

得到如下图像：



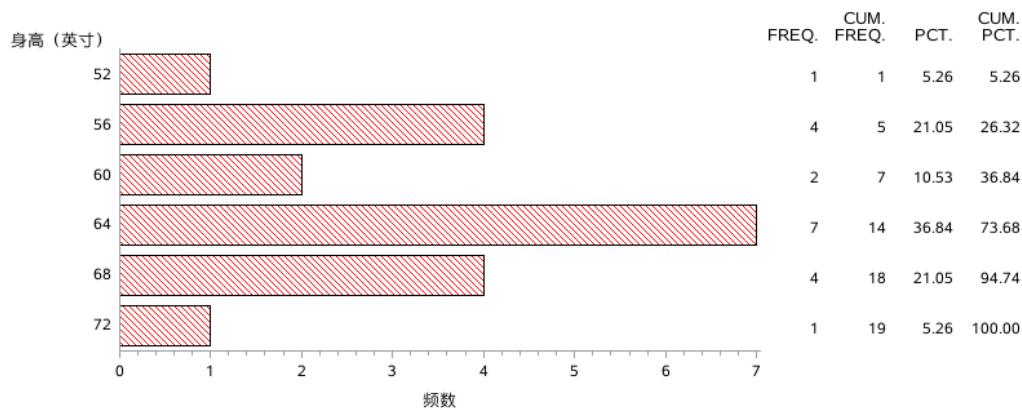
3.9 画 $height$ 的直方图、年龄的饼图

3.9.1 $height$ 的直方图

输入以下SAS代码：

```
1 GOPTIONS RESET = ALL;
2 proc gchart data = sashelp.class;
3     HBAR height / levels = 6 CFRAME = WHITE NOFRAME;
4     PATTERN V = L2, C = red;
5 run;
```

提交后得到如下直方图：



3.9.2 age的饼图

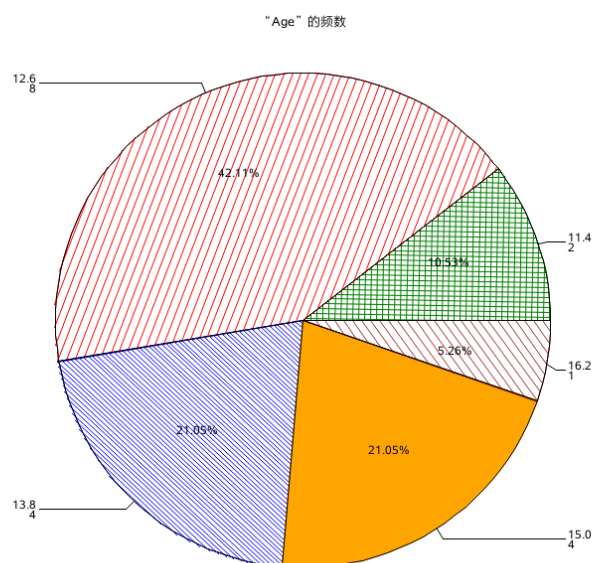
输入以下SAS代码：

```

1  GOPTIONS RESET = ALL;
2  proc gchart data = sashelp.class;
3      PIE age / slice = arrow percent = inside value = arrow;
4      PATTERN1 v = p3x70, c = green;
5      PATTERN2 v = p2n45, c = red;
6      PATTERN3 v = p5, c = blue;
7      PATTERN4 v = ps, c = orange;
8      PATTERN5 v = p2n45, c = brown;
9  run;

```

提交后得到如下饼图：



3.10 按性别分组，计算不同年龄的平均身高

输入以下SAS代码：

```
1 PROC MEANS data = sashelp.class mean ;
2     class age sex;
3     var height;
4 RUN;
```

提交后得到如下表格，即按性别分组后，不同年龄的平均身高：

SAS 系统			
MEANS PROCEDURE			
分析变量: Height 身高			
年龄	性别	观测的个数	均值
11	女	1	51.3000000
	男	1	57.5000000
12	女	2	58.0500000
	男	3	60.3666667
13	女	2	60.9000000
	男	1	62.5000000
14	女	2	63.5500000
	男	2	66.2500000
15	女	2	64.5000000
	男	2	66.7500000
16	男	1	72.0000000

4 反思与总结

4.1 统计分析的几个常用过程

1. *MEANS*: 分类计算变量的常用统计量;
2. *UNIVARIATE*: 计算单变量的统计量和分布的拟合检验（正态检验 histogram /normal lognormal）；
3. *FREQ*: 计算变量取值的频数。

4.2 *FREQ*输出频数表格

```

1 PROC FREQ DATA = 数据集名;
2     TABLES 变量*变量 变量*变量 . . ./
3         nocol (不显示列统计) norow (不显示行) nocum nofreq
4
5         nopercnt missing list
6         out = 数据集 outpct ;
7     WGTGHT 变量名; #按照某个变量加权
8     BY 变量名; #按照某个分类变量统计
9
10 RUN;
```

4.3 定义输出格式的*FORMAT*过程（对连续性变量分组）

```

1 PROC FORMAT library=saslib;
2     VALUE 格式名
3         范围1 = 格式化值1 ...
4         范围n = 格式化值n ;
5 RUN;
```

注1: 对字符型数据定义格式，格式名要加前缀\$。

注2: 若范围超过一个值，离散值可以用，或-用于连续区间，此外<表示到某个临界值范围但不包括该值。

例：

```

'A' = 'Asia'
1,3,5,7,9 = 'Odd'
50000 - HIGH = 'Not Affordable'
13 -< 20 = 'Teenager'
0 <- HIGH = 'Positive Non Zero'
OTHER = 'Bad Data'
```

4.4 *MEANS*—描述统计分析

4.4.1 *MEANS*过程步

```

1  PROC MEANS  DATA = 数据集名 maxdec = 位数 fw = 域宽 noprint 输出统计量
   名列;
2      VAR  变量名列;
3      CLASS  变量名列;
4      BY  变量名列;
5      ID  变量名;
6      OUTPUT  OUT= 数据集名  记入数据集统计量名列;
7  RUN;

```

4.4.2 MEANS过程常用选项:

```

1  MAX      the maximum value
2  MIN      the minimum value
3  MEAN     the mean
4  MEDIAN   the median
5  MODE     the mode(new in sas9.2)
6  N        number of non-missing values
7  NMISS    number of missing values
8  RANGE    the range
9  STDDEV   the standard deviation
10 SUM      the sum
11 CLM      two-sided confidence limits
12 UCLM     upper confidence limits
13 LCLM     lower confidence limits
14 CSS      corrected sum of squares 离差平方和
15 USS      uncorrected sum of squares 平方和
16 ALPHA    confidence level
17 CV       coefficient of variation
18 SKEWNESS  skewness
19 KURTOSIS  kurtosis
20 Q1(P25)   25% quantile
21 Q3(P75)   75% quantile
22 P1,P5,P10,P90,P95,P99
23 STDERR    standard error of mean
24 VAR       variance
25 PROBT     probability for Student's t

```

```

26 T          Student's t
27 MAXDEC     number of   decimal places
28 FW        field width

```

4.5 *UNIVARIATE*过程——单变量特征的概括描述

```

1  PROC UNIVARIATE DATA = 数据集名 noprint ;
2      VAR 变量名列 ;
3      HISTOGRAM 变量名列 / midpoints=中点列
4      normal(mu=均值 sigma=标准差 图象选项)
5      exp(theta= 阈值 图象选项)
6      lognormal(theta=阈值 图象选项) ... .. ;
7      INSET 统计量关键名= '显示名' 格式 ;
8      QQPLOT 变量名列 / square . . . ; %检验是否正态/直线
9      PROBPLOT 变量名列 / . . . ;
10     BY 变量名列 ;
11     ID 变量名 ;
12 RUN ;

```

4.6 上述三种描述性过程的output区别

1. *FREQ*

在table后（选项）（一般输出频数 频数百分比）。

2. *MEANS*（与var并列）

output out=数据集 需要输出的统计量（不设定的话 默认输出n min max mean 并且在同一列）。

（1）自己选择要输出的统计量后就会输出成不同的列；

（2）第一行会输出一行汇总数据（一般无用）。

3. *UNIVARIATE*（与var并列）

output out=数据集 需要输出的统计量（不设定的话 只输出var的值）。