

统计软件期末复习

一、建立（逻辑）库

1. 建库的原因

- A. **逻辑库**就是存放在同一**文件夹**中的一组 SAS 文件, 每一个逻辑库对应了计算机中的一个文件夹;
- B. SAS 数据文件均放在逻辑库中, 使用 SAS 进行对数据操作分析前要指明它在哪一个逻辑库中;
- C. SAS 的逻辑库分为临时库和永久库两种: 临时库名为 **Work**, 存放在 **Work** 中的 SAS 文件叫临时文件, 这些临时文件当退出 SAS 系统时会被自动删除, 永久库中的文件则不会被删除。

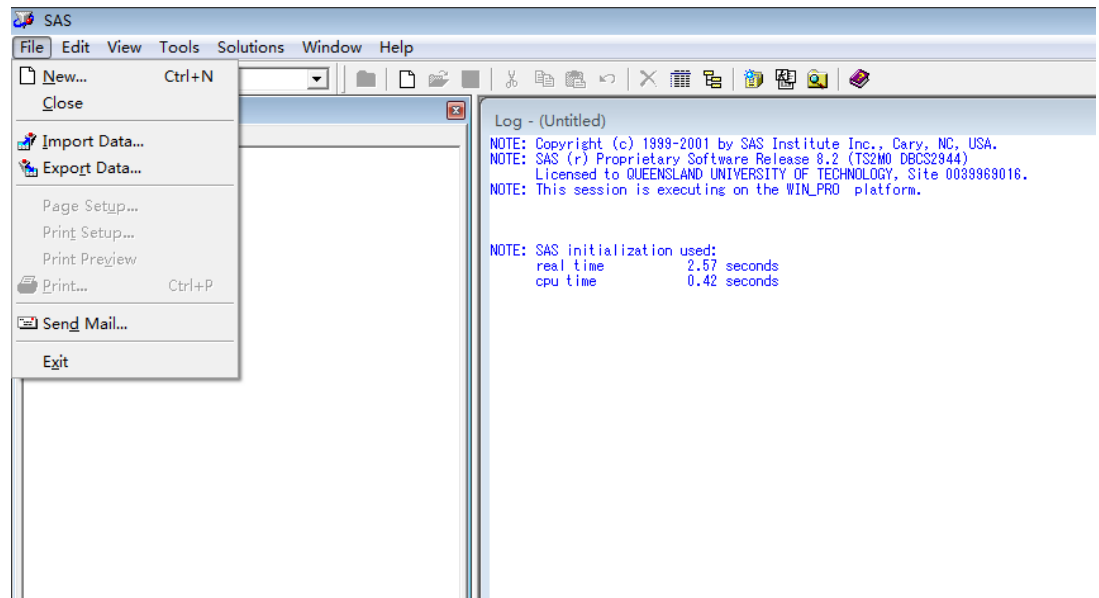
2. 建库操作

1) 编程操作

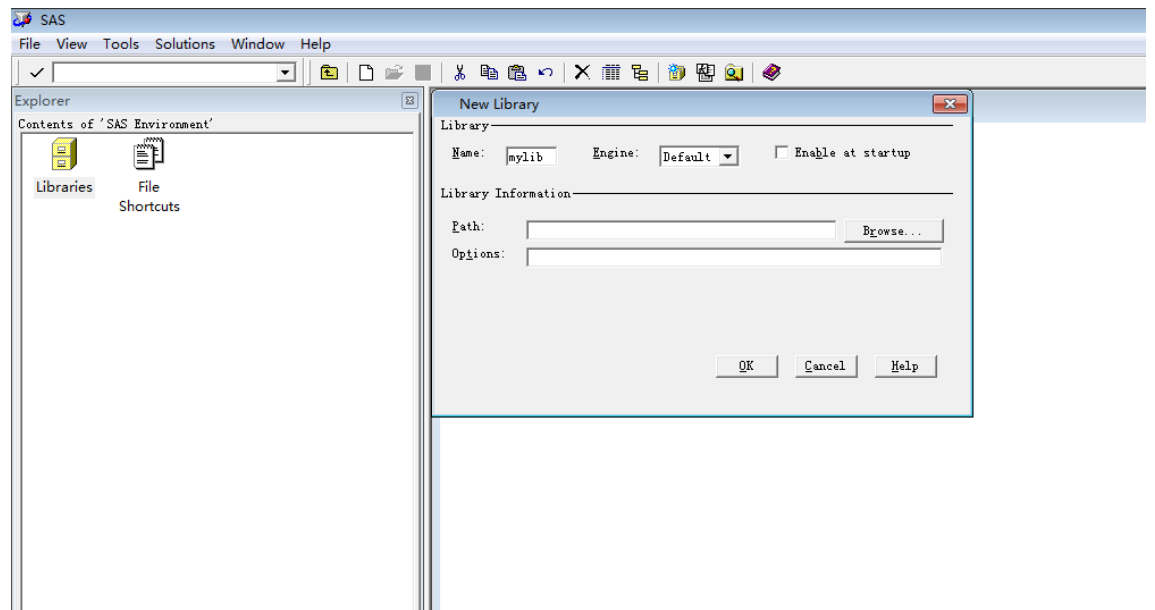
```
libname mysas 'E:\mysas';
```

2) 菜单操作

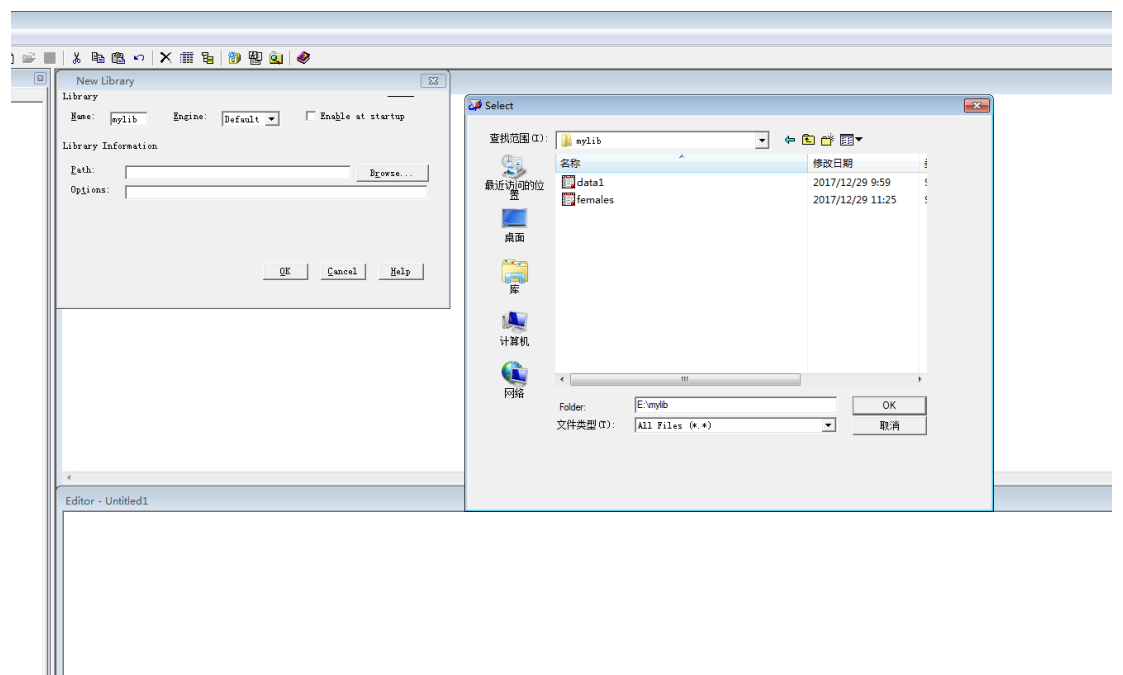
- A. 选择 File 选项的 New 选项



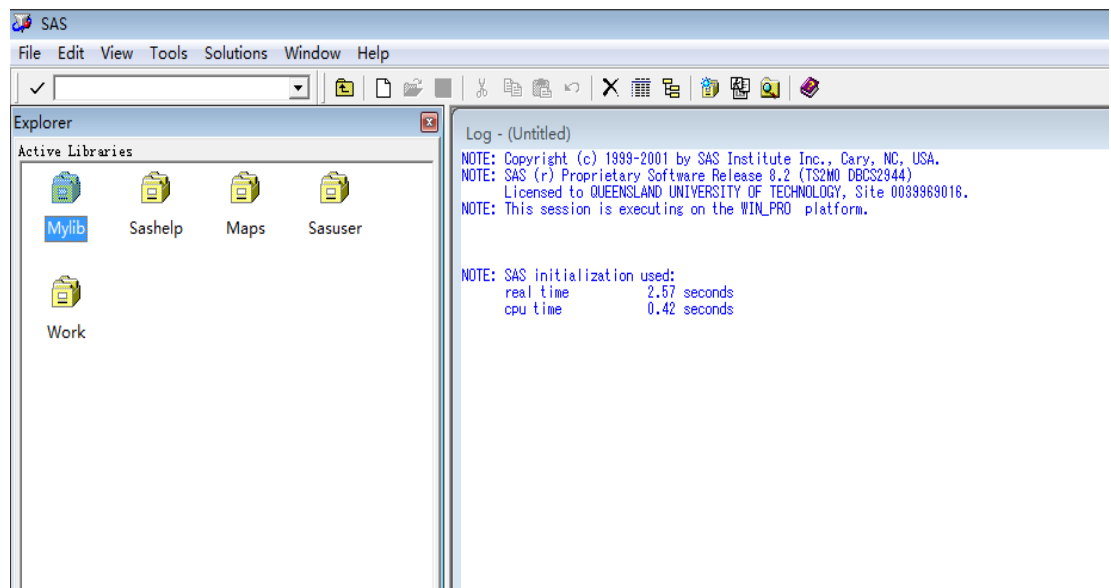
- B. 输入逻辑库名, 选择 **browse** 键



C. 选择对应的文件夹

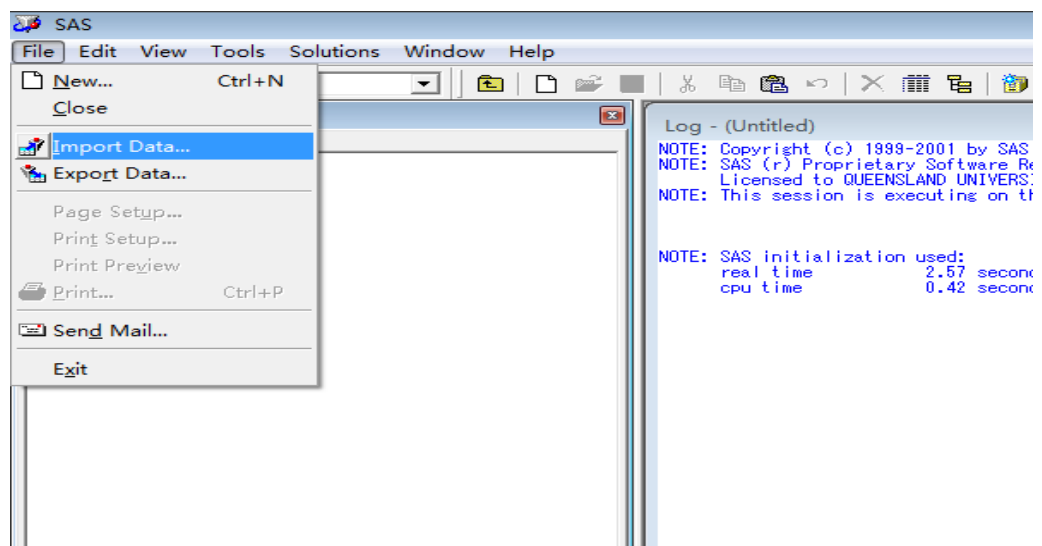


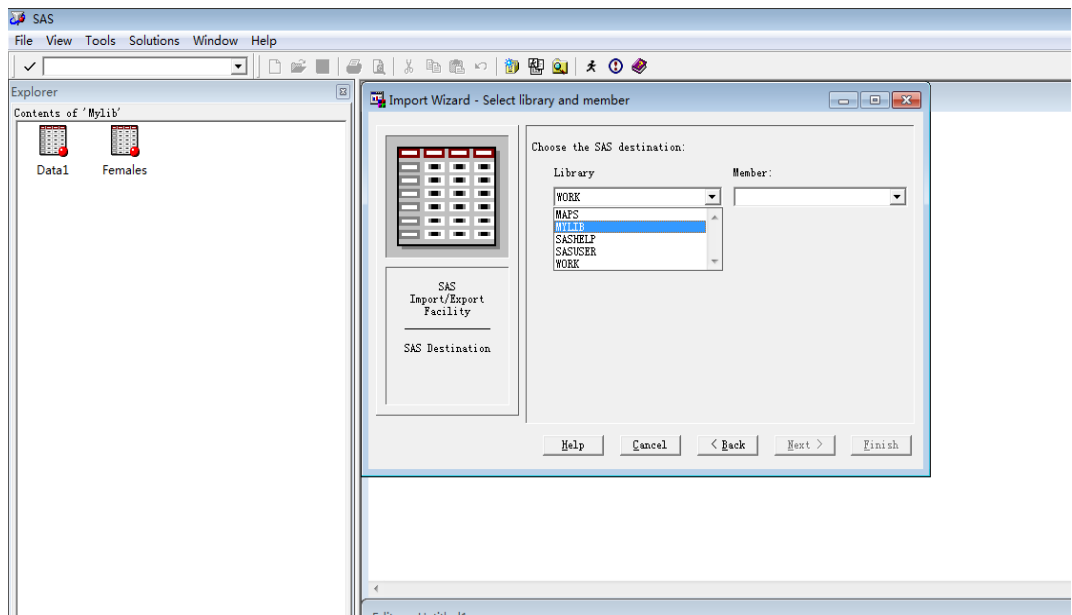
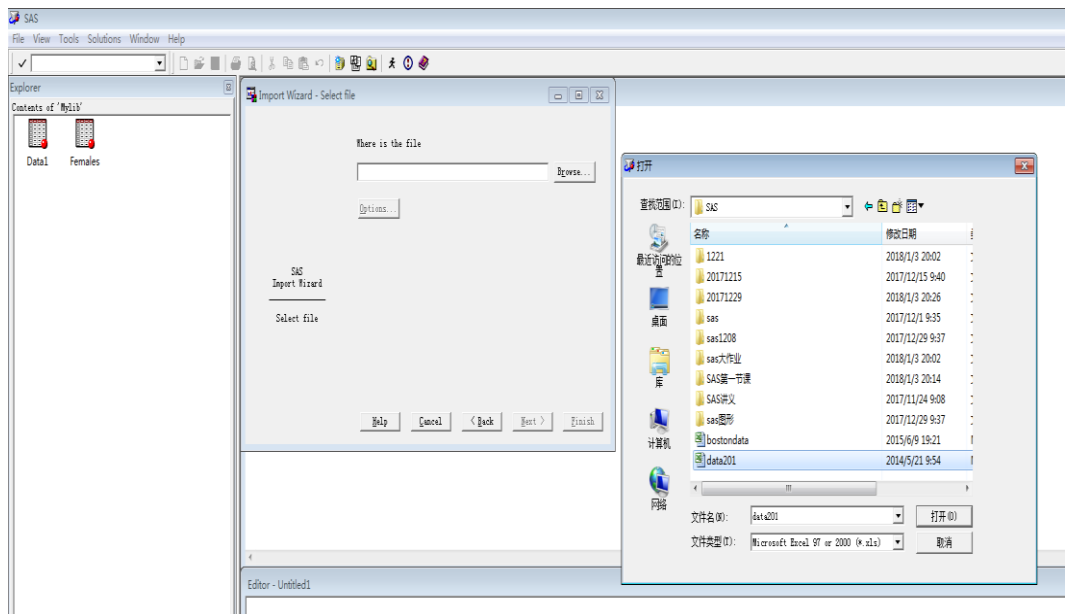
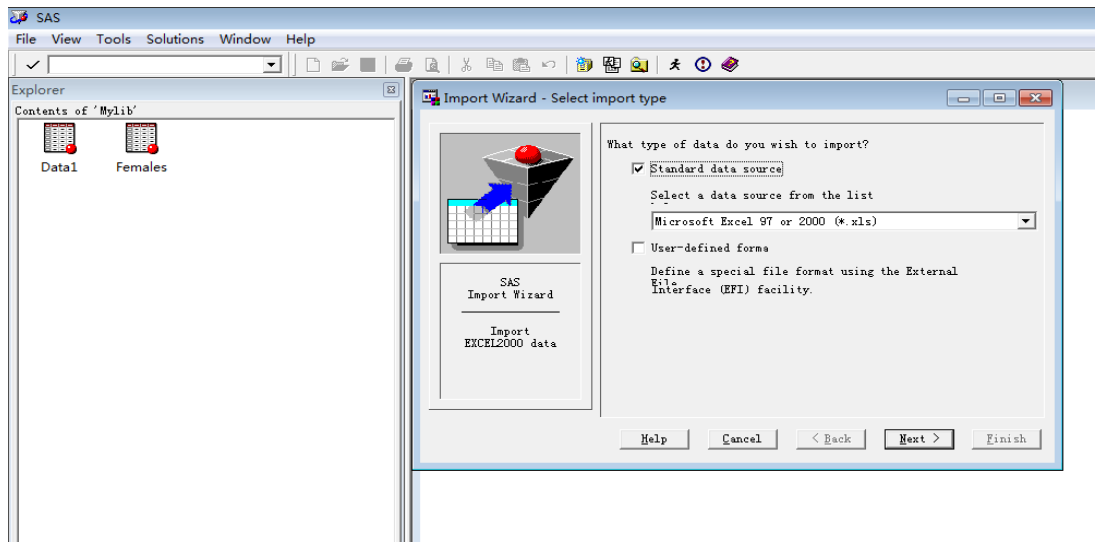
D. 点击 OK，完成建库

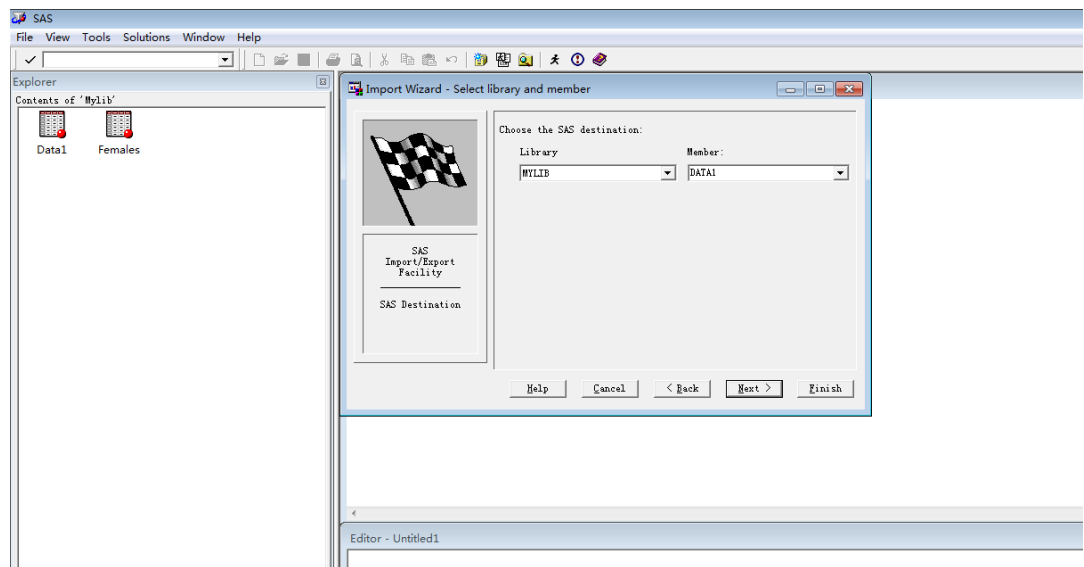


二、导入数据

1. 菜单式操作







2. 编程式操作

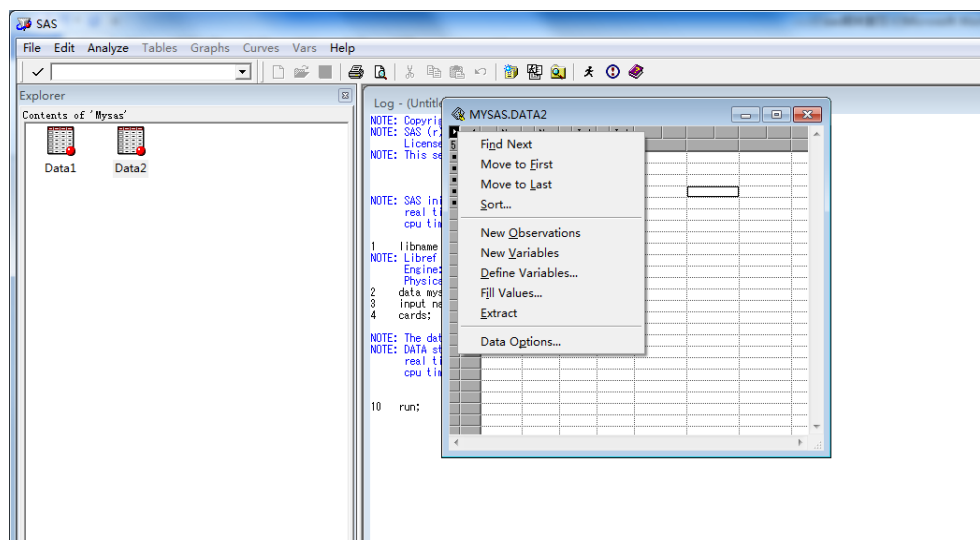
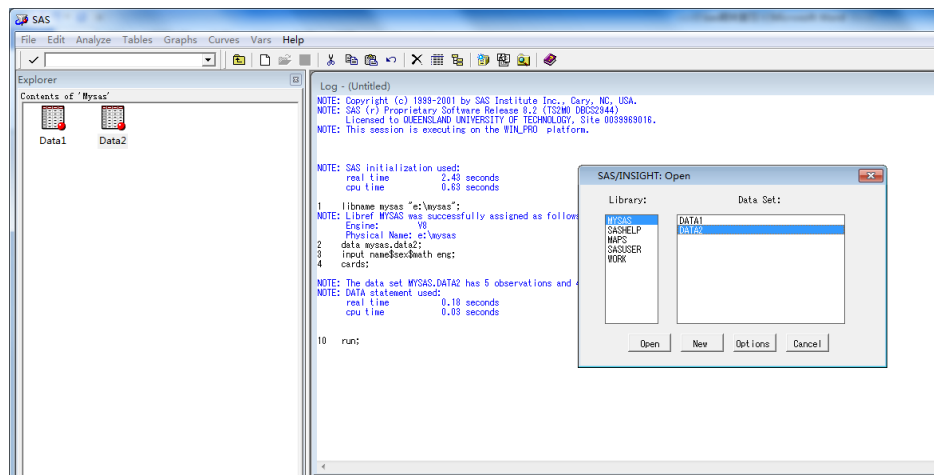
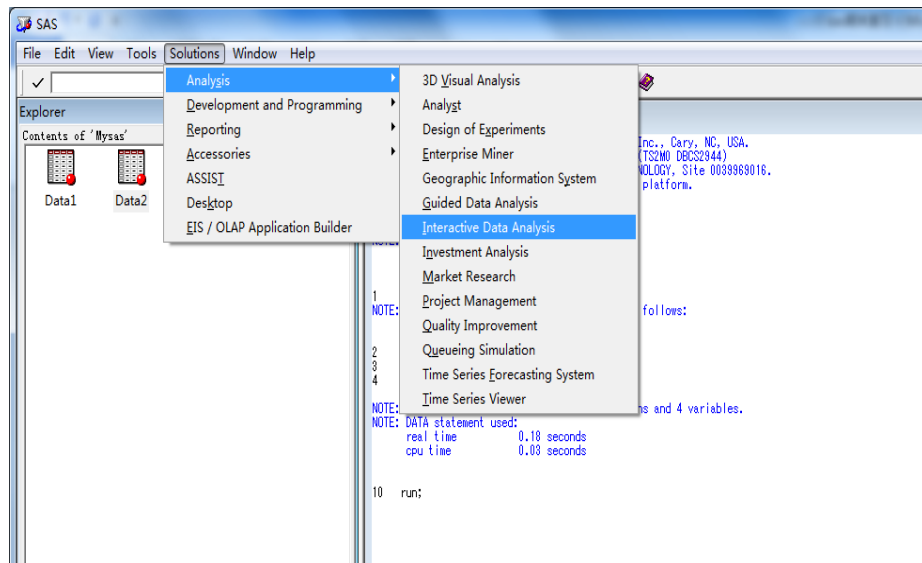
```
proc import datafile="E:\mylib\1.xls" out=mylib.data;
run; /*第一行是变量名*/
```

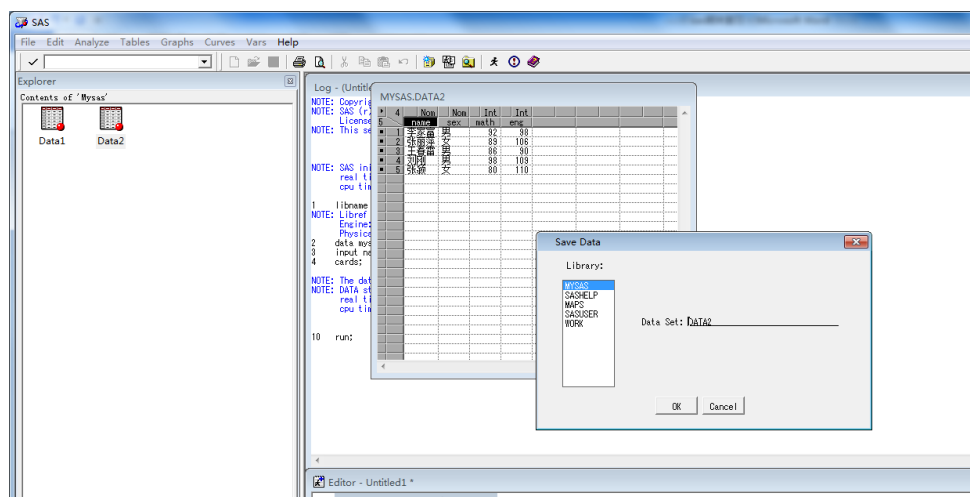
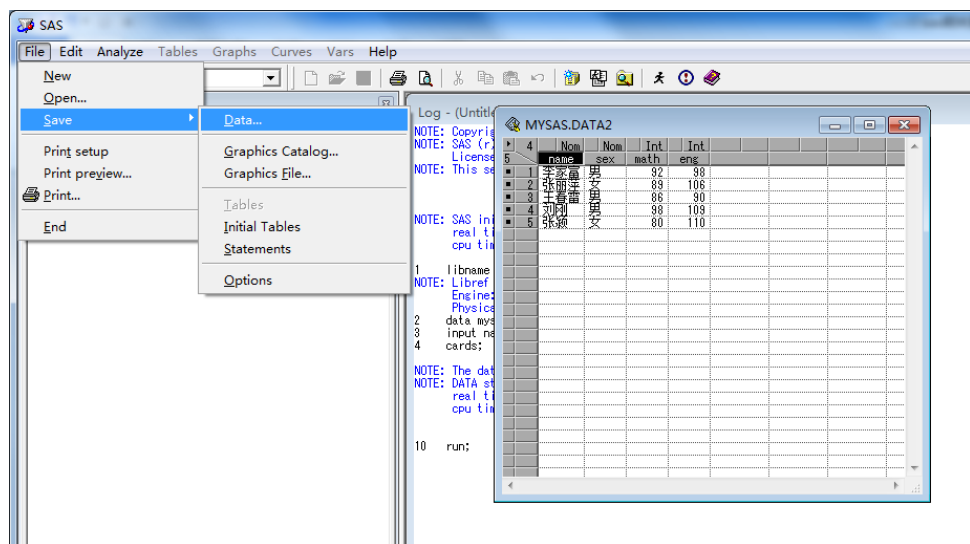
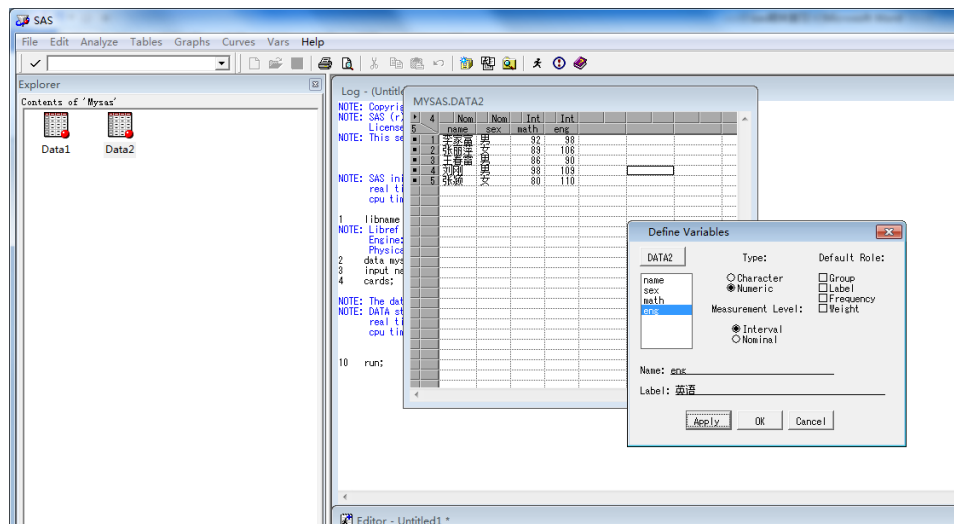
```
libname mysas "e:\mysas";
data mysas.data2;
input name$sex$math eng;
cards;
李家富 男 92 98
张丽萍 女 89 106
王春雷 男 86 90
刘刚 男 98 109
张颖 女 80 110
run;
```

三、对数据的操作

1. 数据的修改

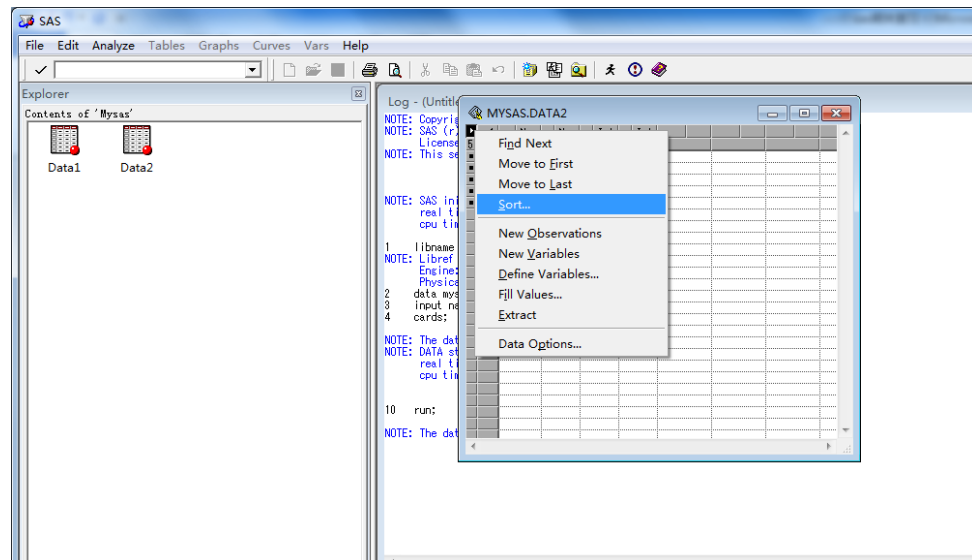
1) 对变量名，变量标签等的修改





注意：修改变量名标签后一定要重新保存数据！（需把观看模式的数据页面关闭才能保存）

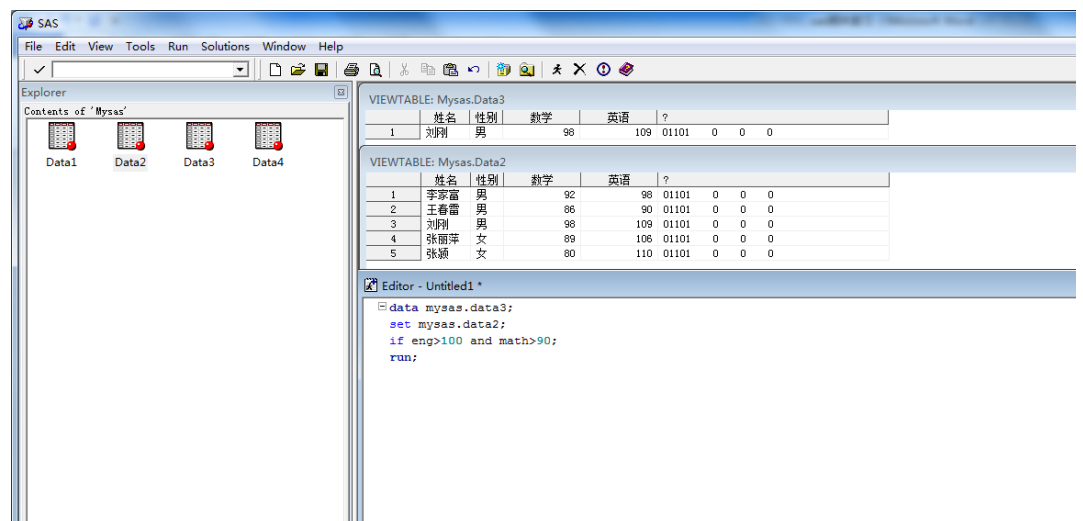
2) 对数据排序



```
proc sort data=mysas.data2;
by sex;
run;
```

3) 筛选数据

```
data mysas.data3;
set mysas.data2;
if eng>100 and math>90;
run;
```



4) 数据的删除, 创建新变量

```
data mysas.data2;
set mysas.data1;
money=income-expense; /*新变量*/
drop num; /*去掉变量*/
run;
```


VIEWTABLE: Mysas.Data1

	家庭编号	地区编号	家庭总收入	家庭总支出	_OBSTAT_			
1	1	2	1794	1550	01101	0	0	0
2	2	2	1761	1365	01101	0	0	0
3	3	1	3410	2730	01101	0	0	0
4	4	2	1765	1530	01101	0	0	0
5	5	2	2184	1900	01101	0	0	0
6	6	2	2050	2050	01101	0	0	0
7	7	2	2460	2184	01101	0	0	0
8	8	1	1976	1170	01101	0	0	0
9	9	1	2850	2496	01101	0	0	0
10	10	1	4275	2760	01101	0	0	0
11	11	2	2010	1275	01101	0	0	0
12	12	1	2236	1810	01101	0	0	0
13	13	1	3305	2820	01101	0	0	0
14	14	1	2400	1976	01101	0	0	0
15	15	2	2250	1970	01101	0	0	0
16	16	2	2200	2060	01101	0	0	0
17	17	1	2730	2236	01101	0	0	0
18	18	1	2496	1455	01101	0	0	0
19	19	1	1760	1040	01101	0	0	0
20	20	1	2820	2366	01101	0	0	0
21	21	2	2250	1966	01101	0	0	0
22	22	1	3170	2400	01101	0	0	0
23	23	2	1200	1250	01101	0	0	0
24	24	?	1776	1350	01101	0	0	0

	地区编号	家庭总收入	家庭总支出	_OBSTAT_	money
1	2	1794	1550	0 0 0	244
2	2	1761	1365	0 0 0	396
3	1	3410	2730	0 0 0	680
4	2	1765	1530	0 0 0	235
5	2	2184	1900	0 0 0	284
6	2	2050	2050	0 0 0	0
7	2	2460	2184	0 0 0	276
8	1	1976	1170	0 0 0	806
9	1	2850	2496	0 0 0	354
10	1	4275	2760	0 0 0	1515
11	2	2010	1275	0 0 0	735
12	1	2236	1810	0 0 0	426
13	1	3305	2820	0 0 0	485
14	1	2400	1976	0 0 0	424
15	2	2250	1970	0 0 0	280
16	2	2200	2060	0 0 0	140
17	1	2730	2236	0 0 0	494
18	1	2496	1455	0 0 0	1041
19	1	1760	1040	0 0 0	720
20	1	2820	2366	0 0 0	454
21	2	2250	1966	0 0 0	284
22	1	3170	2400	0 0 0	770
23	2	1200	1250	0 0 0	-50
24	?	1776	1350	0 0 0	426

Editor - Untitled1 *

```

data mysas.data2;
set mysas.data1;
money=income-expense;
drop num;
run;

```

5) 数据的拆分，合并

```

data mysas.data3 mysas.data4;
set mysas.data2;
select;
when(money > 0) output mysas.data3;
when(money <= 0) output mysas.data4;
otherwise put money='error';
end;
run;

```

VIEWTABLE: Mysas.Data3

	地区编号	家庭总收入	家庭总支出	_OBSTAT_	money
1	2	1794	1550	01101	0 0 0 244
2	2	1761	1365	01101	0 0 0 396
3	1	3410	2730	01101	0 0 0 660
4	2	1765	1530	01101	0 0 0 235
5	2	2184	1900	01101	0 0 0 284
6	2	2460	2184	01101	0 0 0 276
7	1	1976	1170	01101	0 0 0 806
8	1	2850	2406	01101	0 0 0 354

VIEWTABLE: Mysas.Data4

	地区编号	家庭总收入	家庭总支出	_OBSTAT_	money
1	2	2050	2050	01101	0 0 0 0
2	2	1200	1250	01101	0 0 0 -50
3	1	2455	2550	01101	0 0 0 -95
4	2	1080	1380	01101	0 0 0 -300

```

data mysas.data3 mysas.data4;
set mysas.data2;
select;
when(money > 0) output mysas.data3;
when(money <= 0) output mysas.data4;
otherwise put money='error';
end;
run;

```

```

data mysas.data5;
merge mysas.data3 mysas.data3;
    (by x/*根据实际情况，此处填相同的成员名*/)
run;

```

VIEWTABLE: Mysas.Data5

	地区编号	家庭总收入	家庭总支出	_OBSTAT_	money
1	2	1794	1550	01101	0 0 0 244
2	2	1761	1365	01101	0 0 0 396
3	1	3410	2730	01101	0 0 0 660
4	2	1765	1530	01101	0 0 0 235
5	2	2184	1900	01101	0 0 0 284
6	2	2460	2184	01101	0 0 0 276
7	1	1976	1170	01101	0 0 0 806
8	1	2850	2406	01101	0 0 0 354
9	1	4275	2760	01101	0 0 0 1515
10	2	2010	1275	01101	0 0 0 735
11	1	2236	1810	01101	0 0 0 426
12	1	3365	2820	01101	0 0 0 485
13	1	2400	1976	01101	0 0 0 424
14	2	2250	1970	01101	0 0 0 280
15	2	2200	2060	01101	0 0 0 140
16	1	2730	2236	01101	0 0 0 494
17	1	2496	1455	01101	0 0 0 1041
18	1	1760	1040	01101	0 0 0 720
19	1	2820	2366	01101	0 0 0 454
20	2	2250	1966	01101	0 0 0 284
21	1	3170	2400	01101	0 0 0 770
22	2	1776	1350	01101	0 0 0 426
23	2	1980	1794	01101	0 0 0 186
24	?	1466	1700	01101	n n n 766

```

data mysas.data5;
merge mysas.data3 mysas.data3;
run;

```

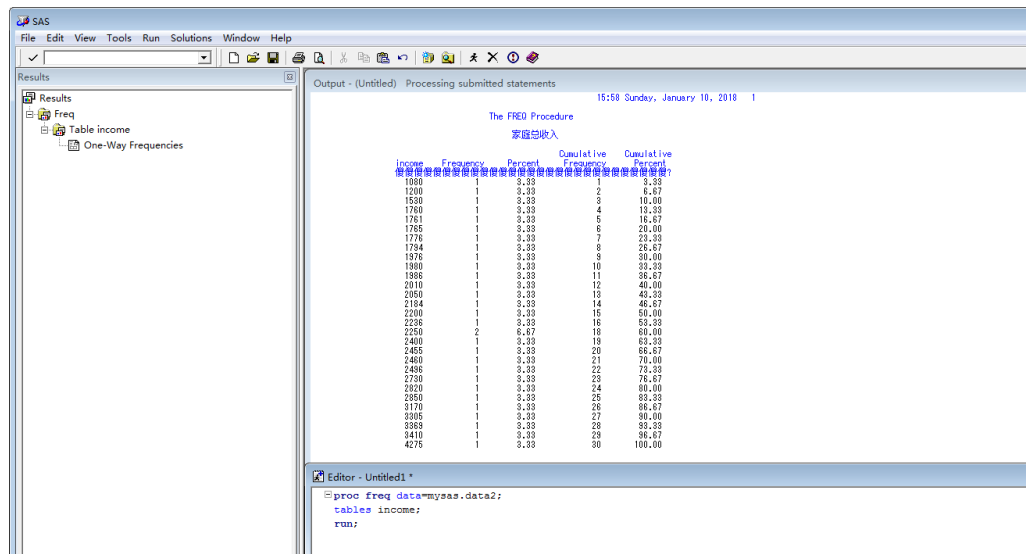
2. 数据的描述统计量
 - 1) 输出描述统计量

Freq 过程

```

proc freq data=mysas.data2;
tables income;
run;

```

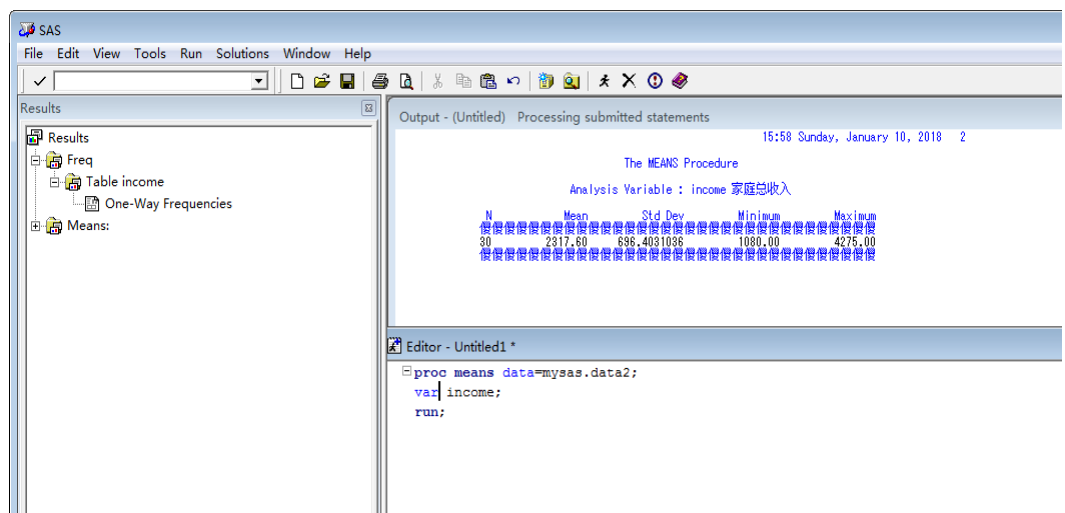


Means 过程

```

proc means data=mysas.data2;
var income; /*输出简单统计量*/
run;

```



```

proc means data=mysas.data2 n mean median p1 p5 p95 p99 q1 q3
max min; /*指定输出统计量，数目，均值，中位数，第一位百分数，第五位百分数，
第九十五位百分数，四分之一分位数，四分之三分位数*/
var income;
run;

```

统计参数的关键词	含 义	统计参数的关键词	含 义
N	样本数	Cv	变异系数
Mean	平均数	Var	方差
Std	标准差	Stderr	均值的标准误
Min	最小值	Skewness	偏度
Max	较大值	Kurtosis	峰度
Nmiss	缺失值个数	Q1 P25	四分之一分位数
Mode	众数	Q3 P75	四分之三分位数

Median	中位数	P1	第 1 百分位数
Range	极差	P5	第 5 百分位数
Uss	加权平方和	P10	第 10 百分位数
Css	均值偏差的 加权平方和	P90	第 90 百分位数
Uclm	置信度上限	P95	第 95 百分位数
Lclm	置信度下限	P99	第 99 百分位数
C1m	置信度上限和下限	QRANGE	百分位数极差
Sum	累加和	PROBT PRT	T 分布的双尾 p 值
Sumwgt	权数和	T	总体均值为 0 的 t 统计量

Univariate 过程

```
proc univariate data=mysas.data2;
var income;
class area;
run;
```

The screenshot displays the SAS Univariate process results. The left pane shows the Results tree with 'Univariate: income' expanded. The right pane shows the Output window with statistical tests and quantiles. A small editor window shows the SAS code used to generate the results.

Tests for Location: Mu0=0

Test	-Statistic-	Pr > t	Pr >= M	Pr >= S
Student's t	t 20.03836	Pr > t	<.0001	<.0001
Sign	M 8	Pr >= M	<.0001	<.0001
Signed Rank	S 68	Pr >= S	<.0001	<.0001

Quantiles (Definition 5)

Quantile	Estimate
100% Max	2460
95%	2460
90%	2460
90% Q3	2250
75% Q3	2192
50% Median	1983

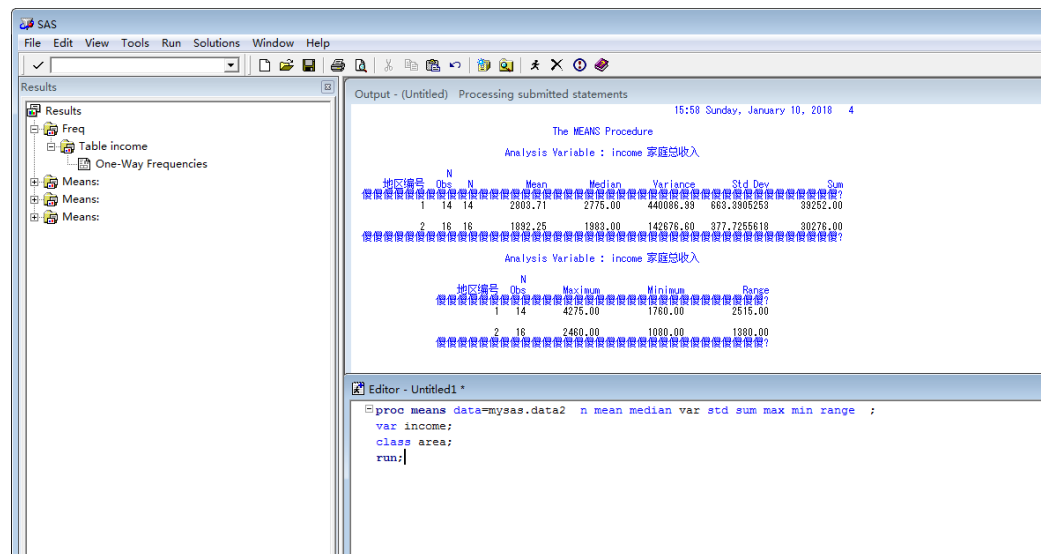
Editor - Untitled1 *

```
proc univariate data=mysas.data2;
var income;
class area;
run;
```

2) 分组输出描述统计量

Class 法

```
proc means data=mysas.data2 n mean median var std sum max min
range ;
var income;
class area;
run;
```

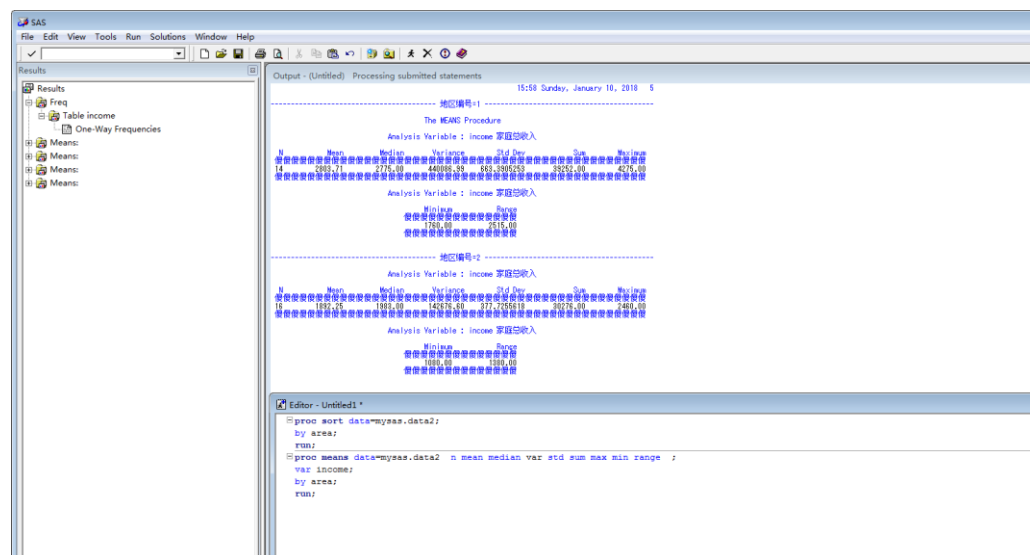


By 法（注意要先排序）

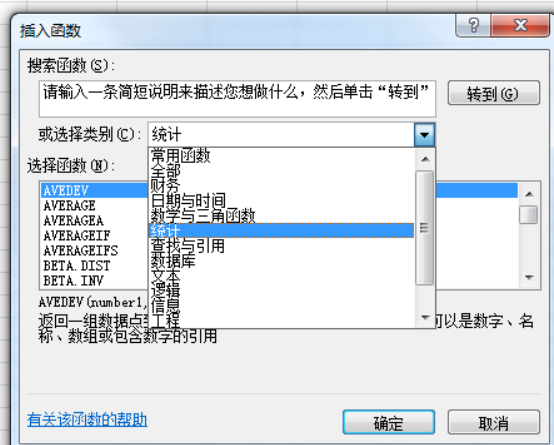
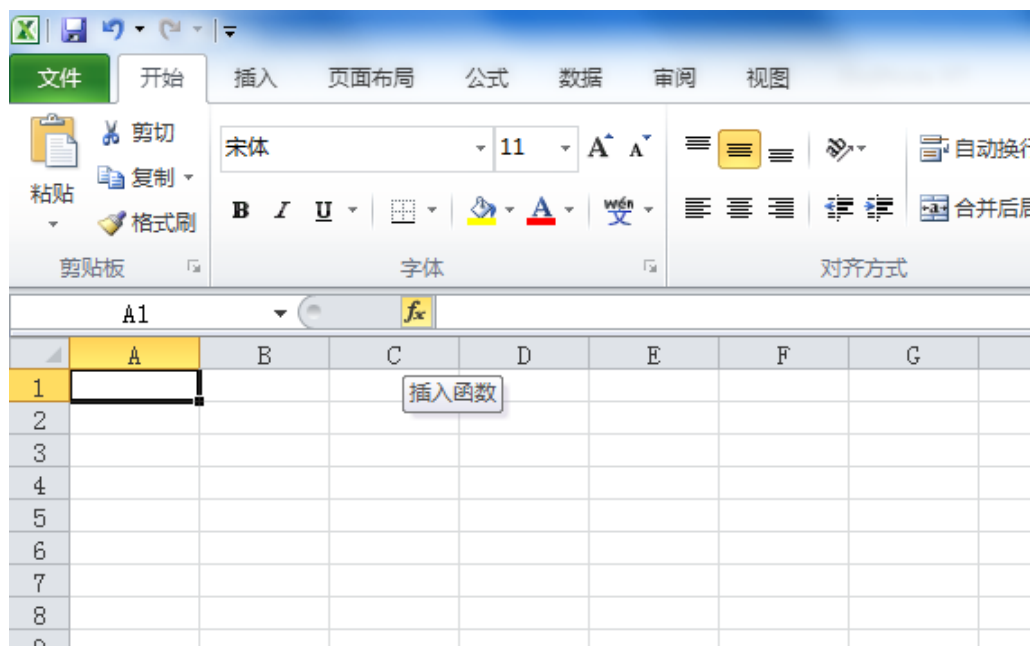
```

proc sort data=mysas.data2;
by area;
run;
proc means data=mysas.data2  n mean median var std sum max min
range ;
var income;
by area;
run;

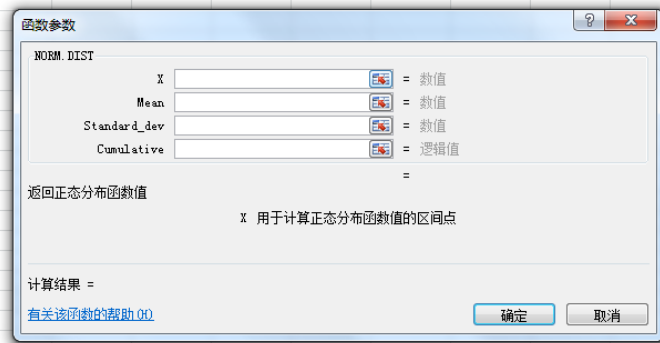
```



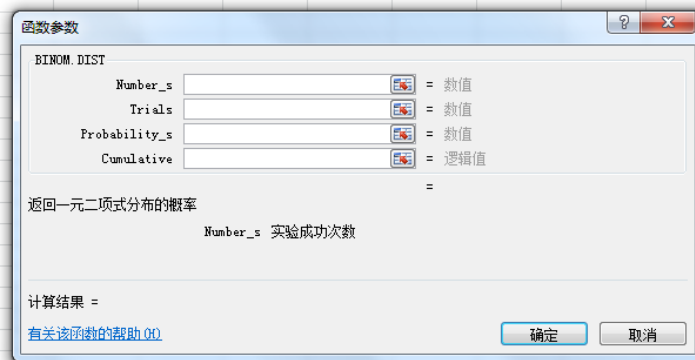
3) 三个典型分布（Excel）



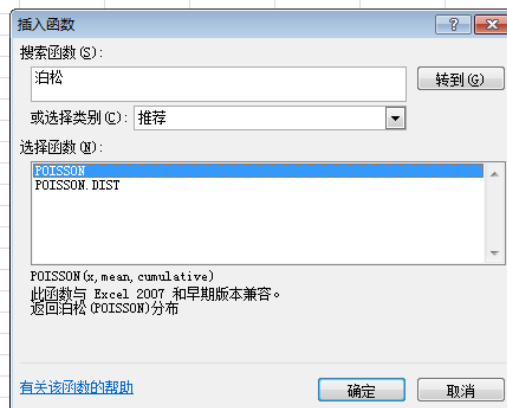
A. 正态分布



B. 二项分布



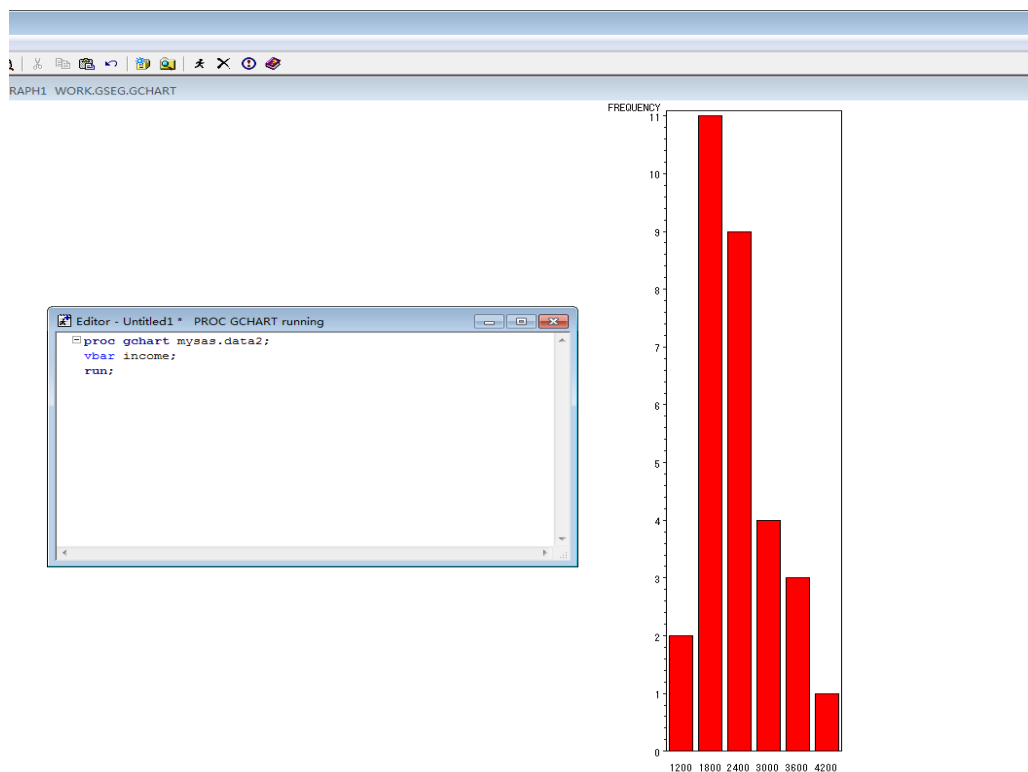
C. 泊松分布



3. 数据的图形描述

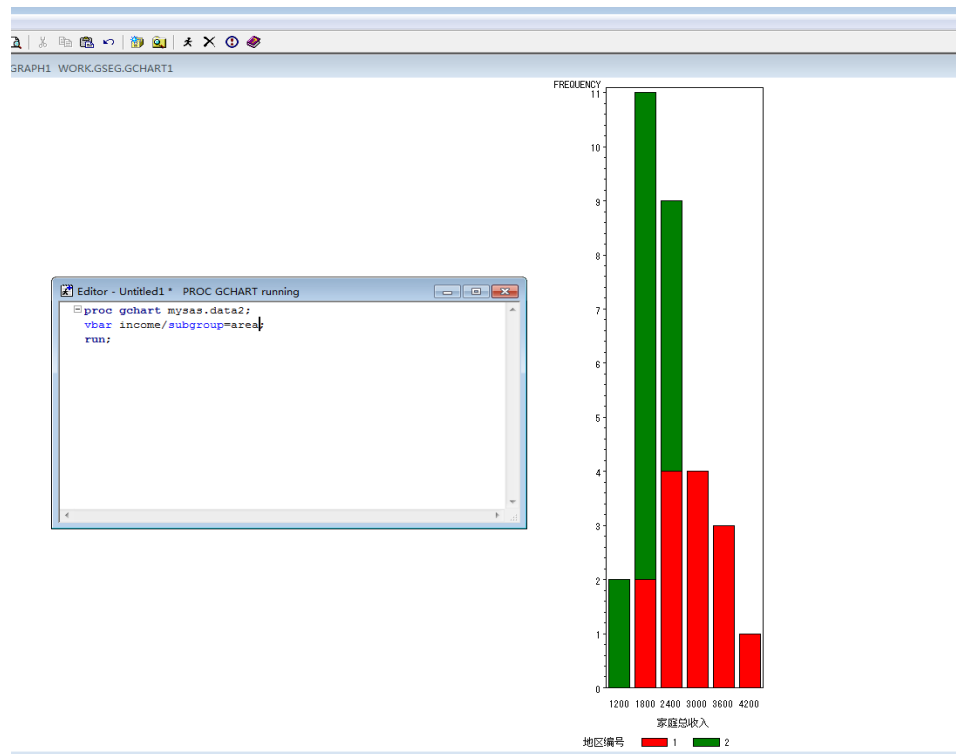
1) 条形图

```
proc gchart mysas.data2;
vbar income;
run;
```



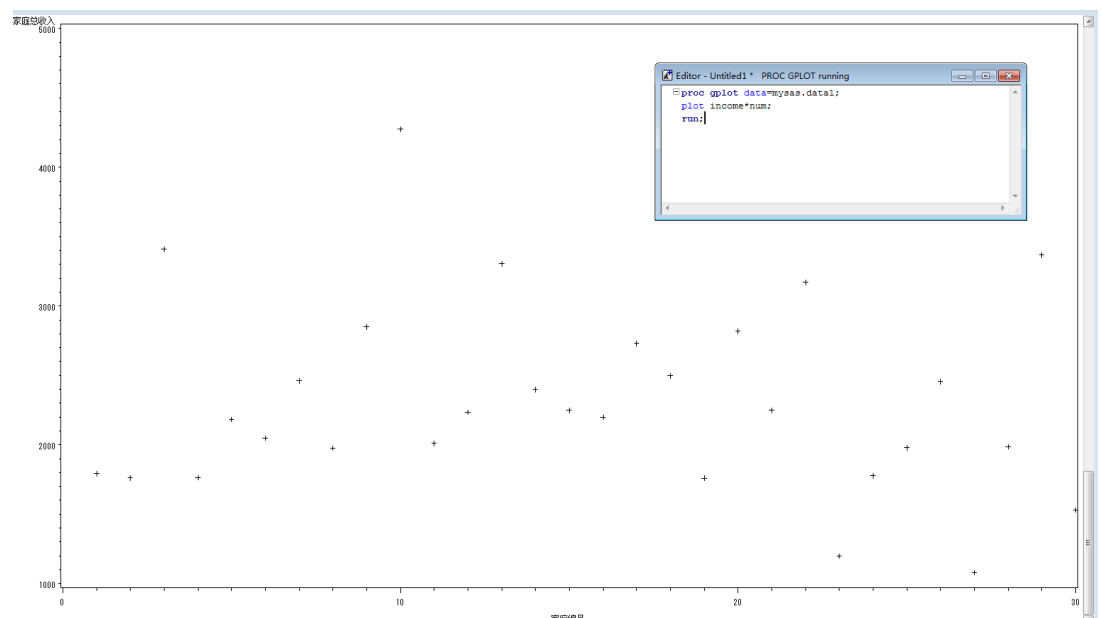
分组条形图

```
proc gchart mysas.data2;
vbar income/subgroup=area;
run;
```

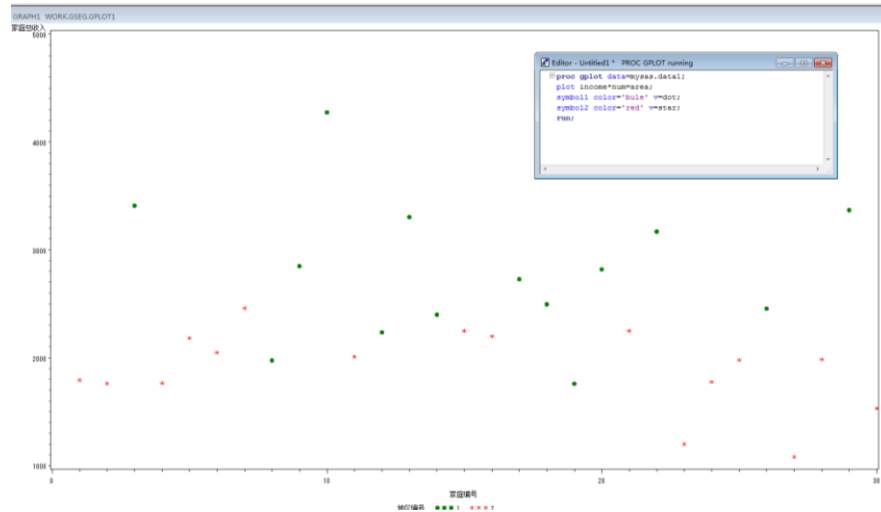



2) 散点图

```
proc gplot data=mysas.data1;
plot income*num; (先y后x)
run;
```



```
proc gplot data=mysas.data1;
plot income*num=area;
symbol1 color='bule' value=dot;
symbol2 color='red' value=star;
run;
```

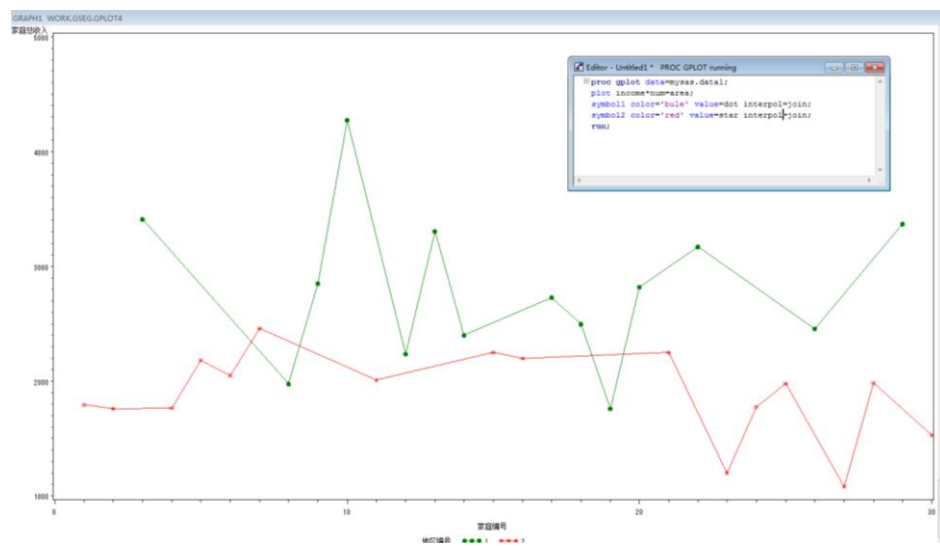


3) 折线图

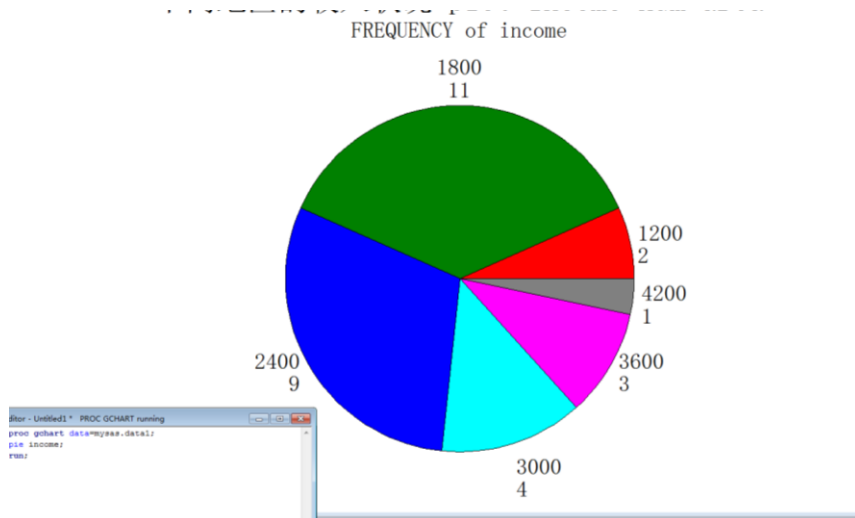
```

proc gplot data=mysas.data1;
plot income*num=area;
symbol1 color='blue' value=dot interpol=join; /*color可用cv代替value可用v代替 Interpol可用i代替*/
symbol2 color='red' value=star interpol=join;
run;

```



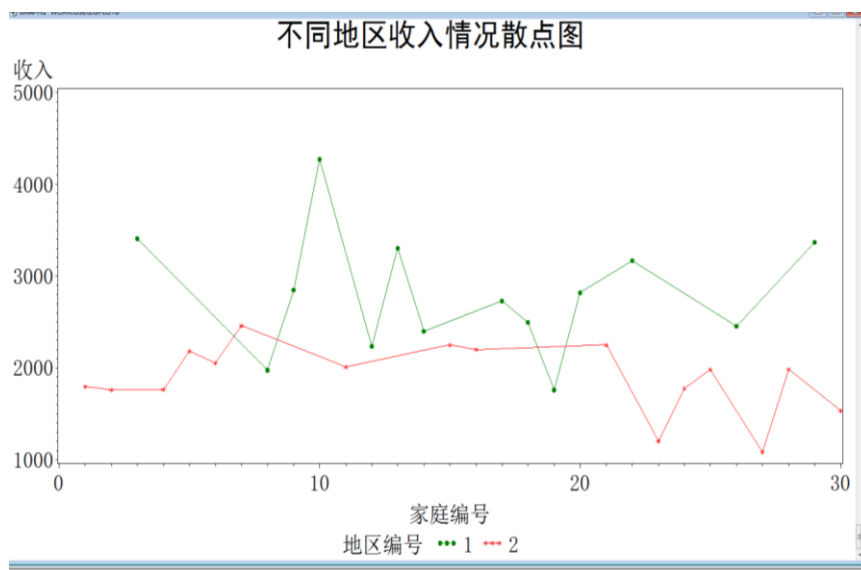
4) 饼图



附：添加标题，修改图形标签等

```

title f='黑体' '不同地区收入情况散点图';
proc gplot data=mysas.data1;
plot income*num=area;
symbol1 color='bule' value=dot;
symbol2 color='red' value=star;
label income='收入' ;
run;
  
```



PS:

整理较匆忙，如果有错误或者遗漏的重要知识点，欢迎补充。

By 凌宛莹

2018.1