



Fudan University Library **2020**
Regular Training

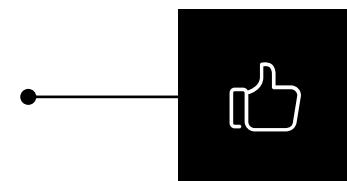
SPSS数据分析基础

■ 胡杰

■ 2020-06-02

如何获取正版SPSS

复旦信息办

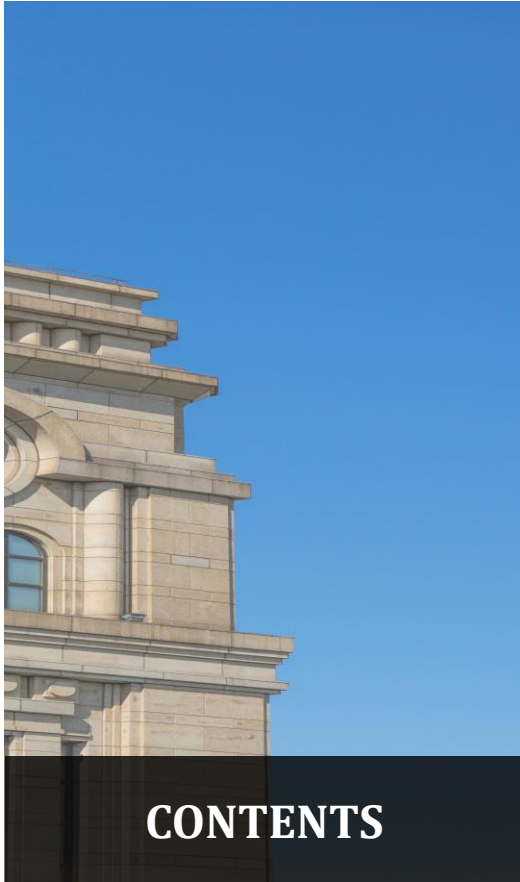


服务对象：复旦大学在校师生

软件下载地址及安装说明请访问：

<http://mvls.fudan.edu.cn>

版本：SPSS 20.0



1 • SPSS快速入门

2 • 数据分析的一般步骤

3 • 问卷数据分析方法



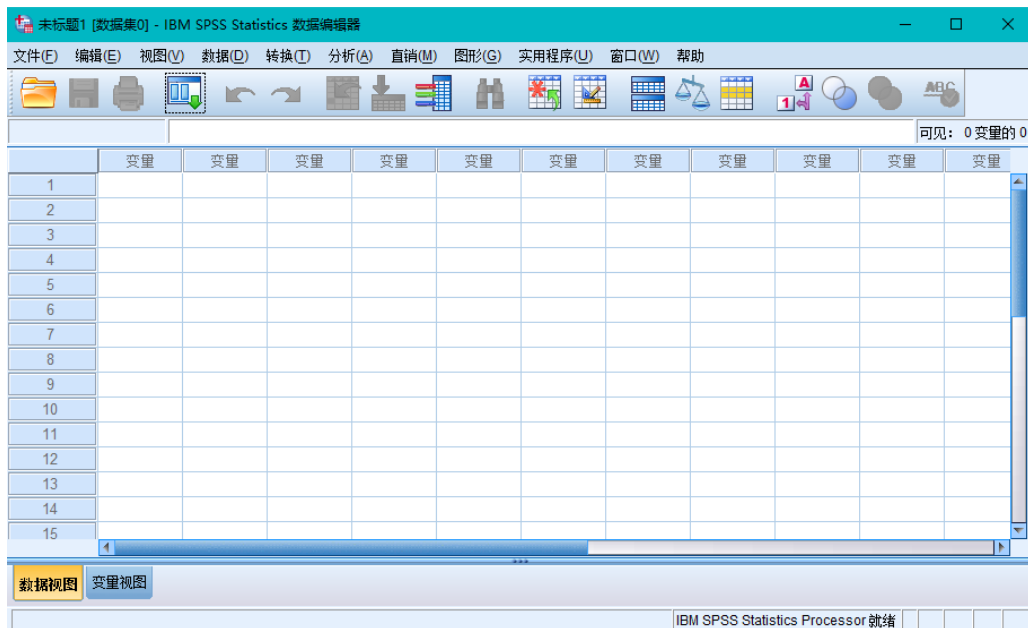
/01

SPSS快速入门

了解SPSS的基本窗口是学习使用SPSS的入门点

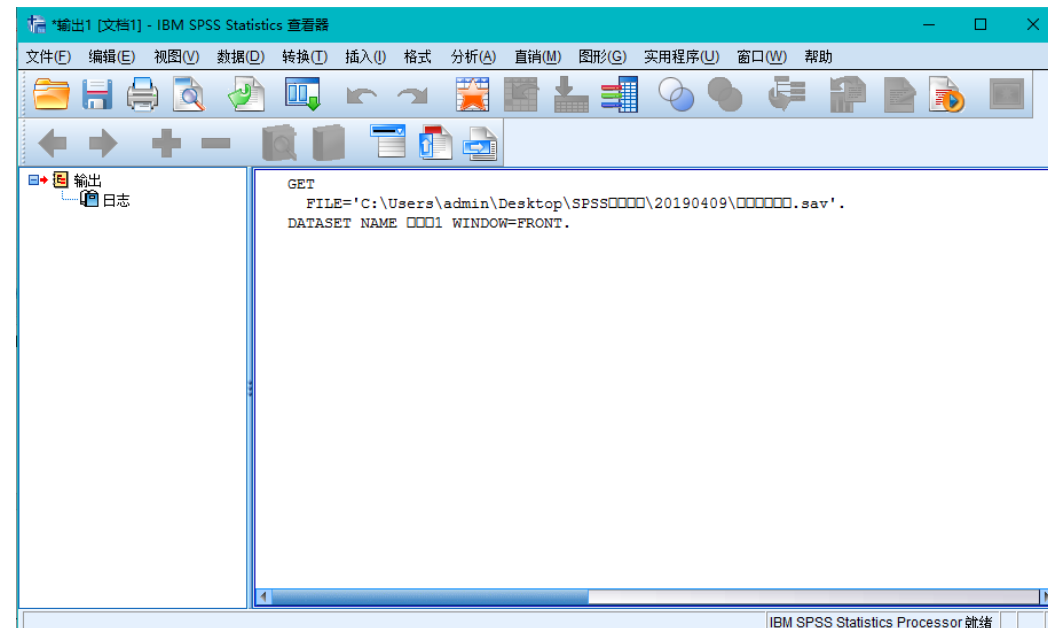
- SPSS有哪些基本操作窗口
- 各个窗口的功能特点是什么
- 窗口之间有怎样的关系

SPSS的基本窗口



数据编辑窗口

(输入和管理待分析数据)



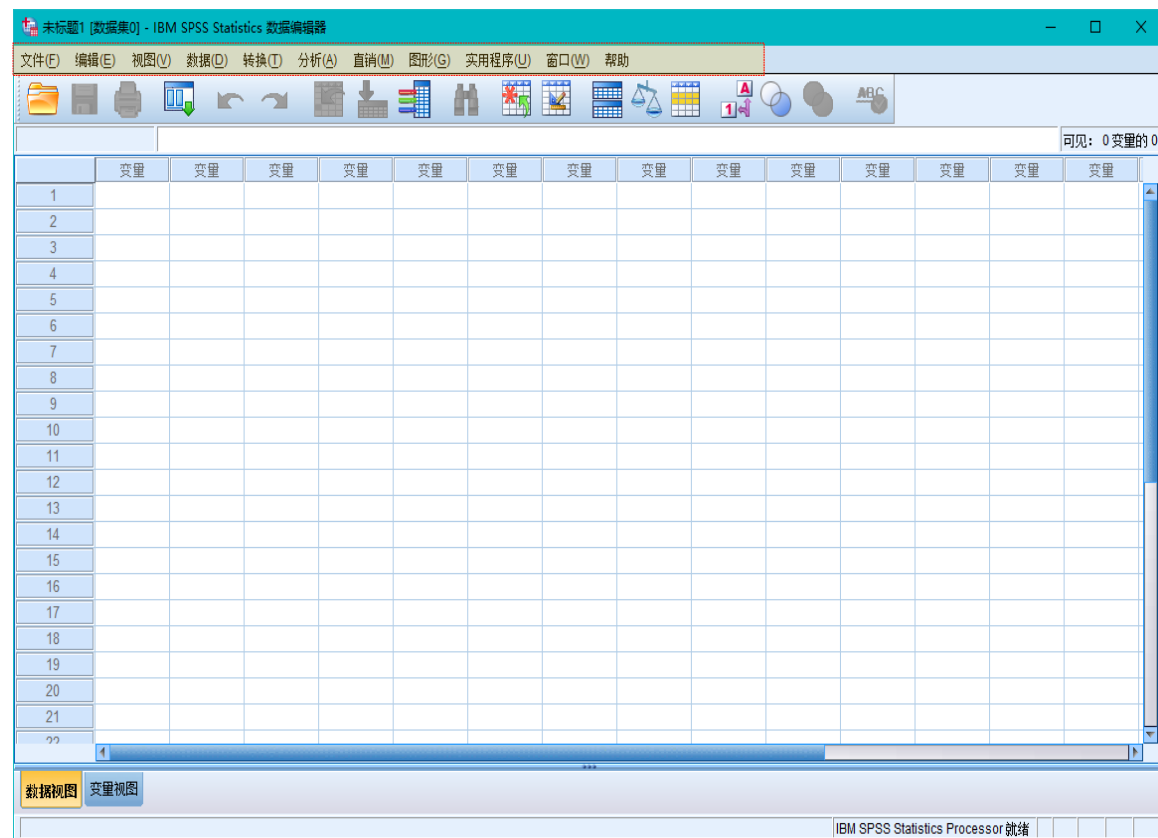
结果输出窗口

(接受和管理统计分析的结果)

数据编辑窗口

窗口主菜单

罗列了SPSS常用的数据编辑、加工和分析的功能



窗口主菜单

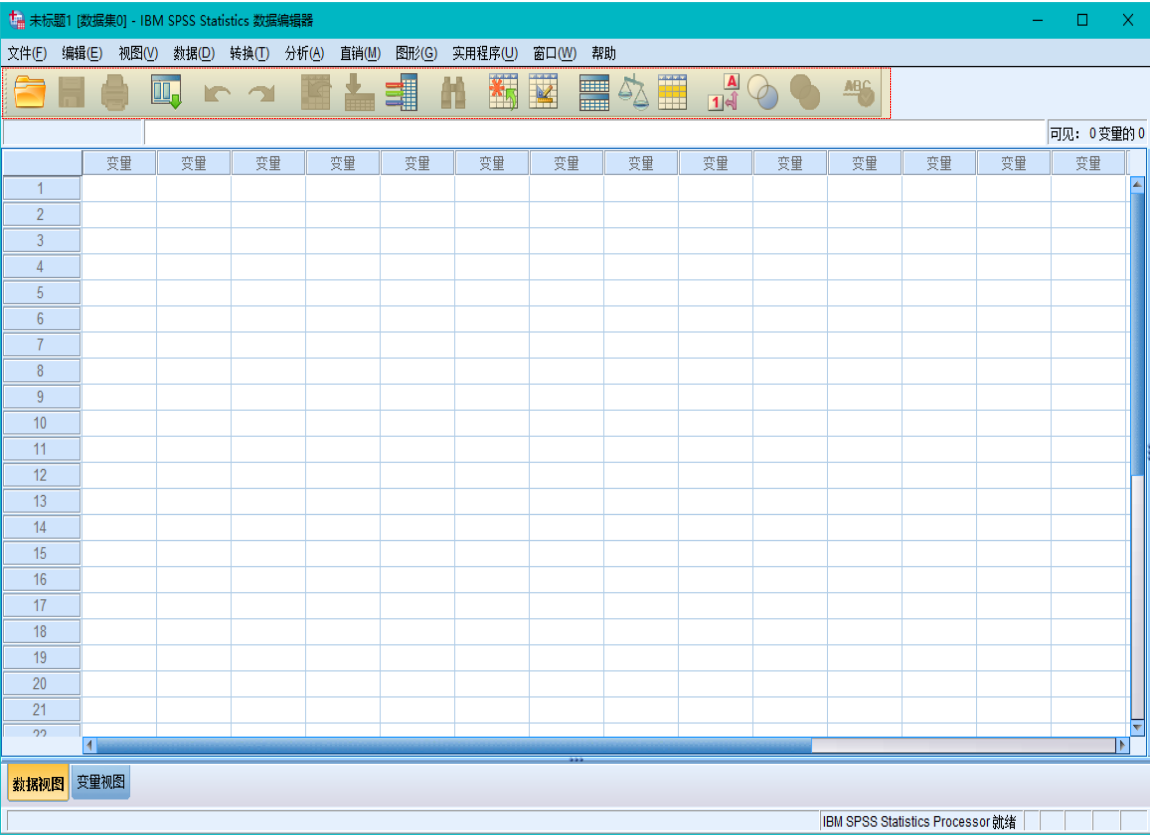


菜单名	功能	解释
文件 (File)	文件操作	对SPSS相关文件进行基本管理：如新建、打开、保存、打印等
编辑 (Edit)	数据编辑	对数据进行基本编辑，实现数据查找等功能：如撤销/恢复、剪切、复制、粘贴等
视图 (View)	窗口外观状态管理	对SPSS窗口外观等进行设置：如状态栏、表格线等是否显示，字体设置等
数据 (Data)	数据的操作和管理	对数据进行加工整理：如数据的排序、转置、抽样、分类汇总等
转换 (Transform)	数据基本处理	对数据进行基本处理：如生成新变量、计数、分组等
分析 (Analyze)	统计分析	对数据进行统计分析和建模：如基本统计分析、均值比较、回归分析等
图形 (Graphs)	制作统计图形	对数据生成各种统计图形：如条形图、直方图、饼图、线图、散点图等
实用程序 (Utilities)	实用程序	SPSS其他辅助管理：如显示变量信息、定义变量集、菜单编辑器等
窗口 (Windows)	窗口管理	对SPSS的多个窗口进行管理：如窗口切换、最小化窗口等
帮助 (Help)	帮助	实现SPSS的联机帮助：如语句检索、统计教程等

数据编辑窗口

窗口主菜单

罗列了SPSS常用的数据编辑、加工和分析的功能



工具栏

一些常用功能以图形按钮的形式组织在工具栏中，操作更加快捷和方便。

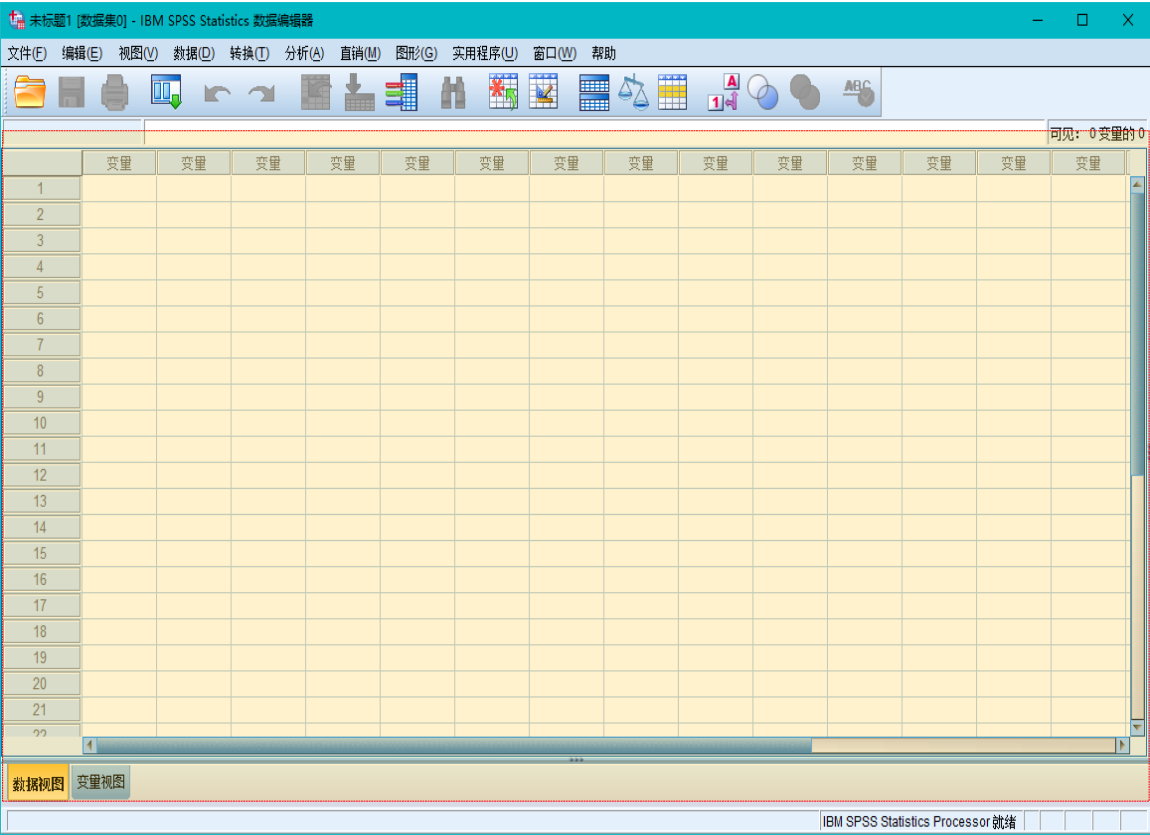
数据编辑窗口

窗口主菜单

罗列了SPSS常用的数据编辑、加工和分析的功能

数据编辑区

数据视图：录入、编辑和管理数据内容
变量视图：定义和修改数据结构



工具栏

一些常用功能以图形按钮的形式组织在工具栏中，操作更加快捷和方便。

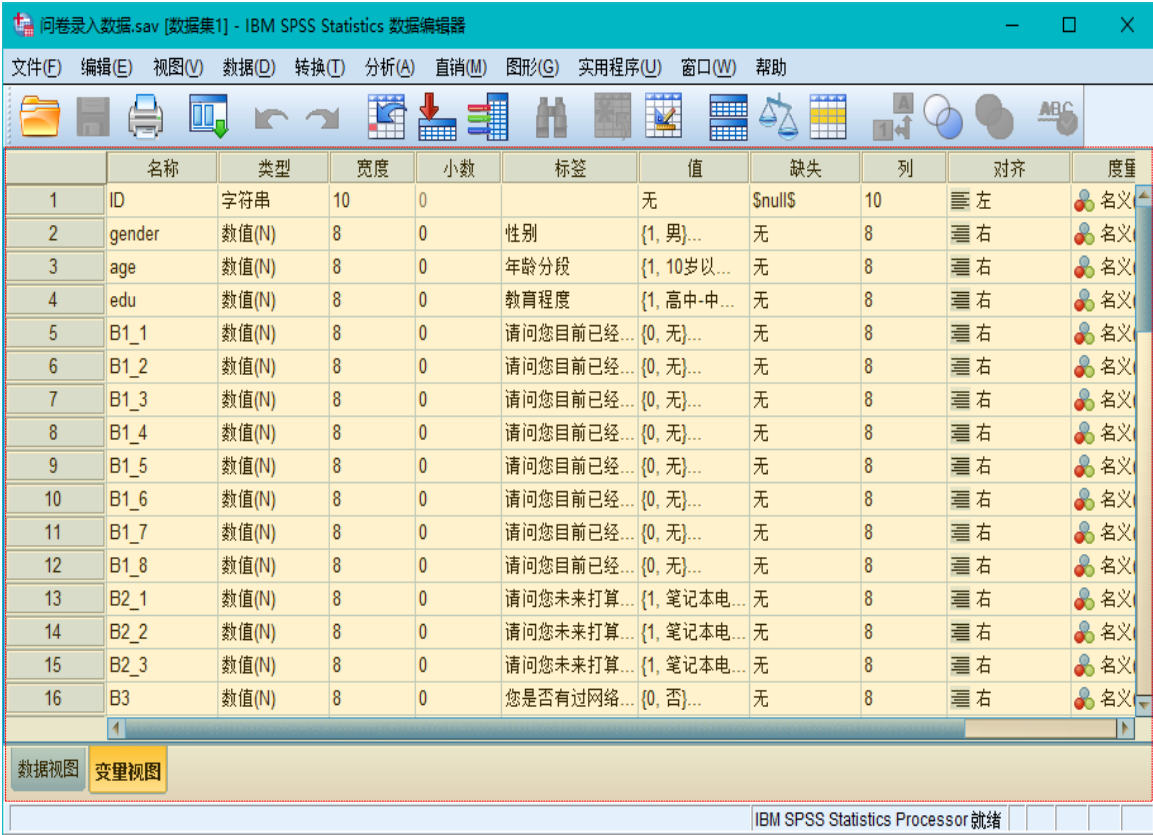
数据编辑窗口

窗口主菜单

罗列了SPSS常用的数据编辑、加工和分析的功能

数据编辑区

数据视图：录入、编辑和管理数据内容
变量视图：定义和修改数据结构



工具栏

一些常用功能以图形按钮的形式组织在工具栏中，操作更加快捷和方便。

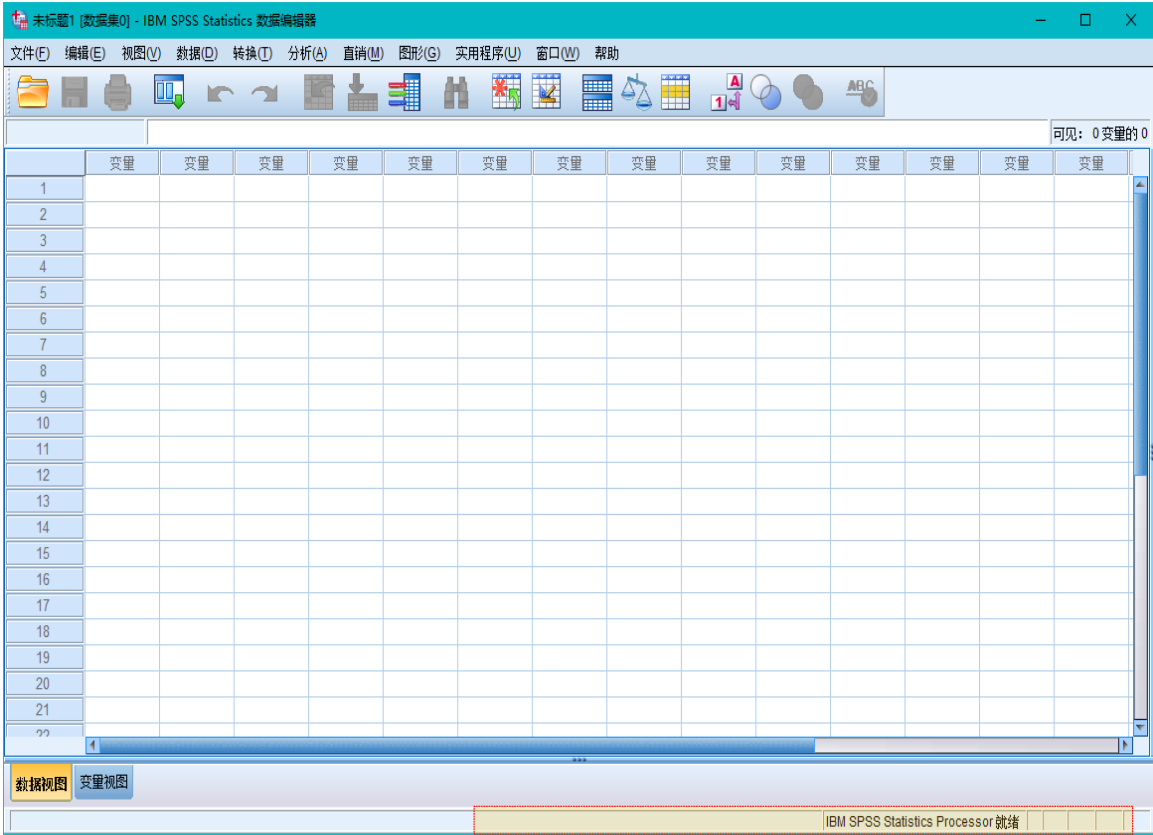
数据编辑窗口

窗口主菜单

罗列了SPSS常用的数据编辑、加工和分析的功能

数据编辑区

数据视图：录入、编辑和管理数据内容
变量视图：定义和修改数据结构



工具栏

一些常用功能以图形按钮的形式组织在工具栏中，操作更加快捷和方便。

系统状态显示区

显示系统的当前运行状态

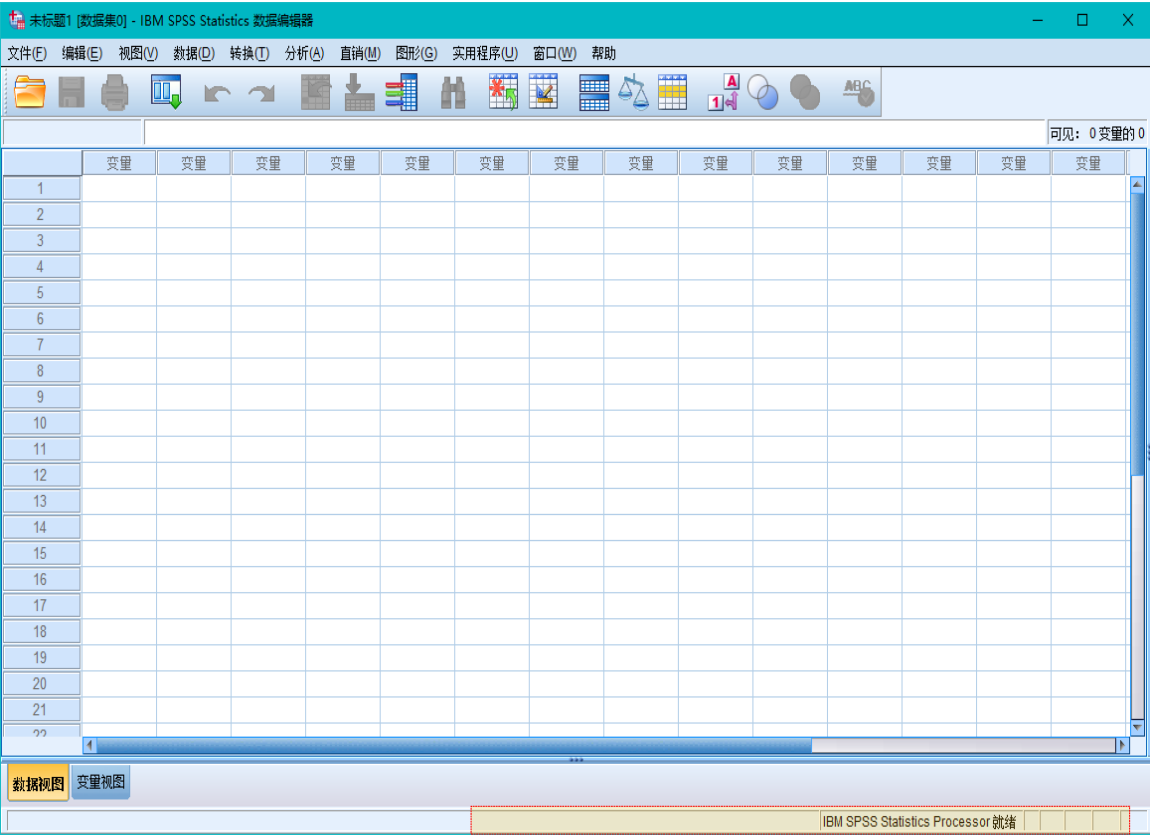
数据编辑窗口

窗口主菜单

罗列了SPSS常用的数据编辑、加工和分析的功能

数据编辑区

数据视图：录入、编辑和管理数据内容
变量视图：定义和修改数据结构



工具栏

一些常用功能以图形按钮的形式组织在工具栏中，操作更加快捷和方便。

系统状态显示区

显示系统的当前运行状态

- ✓ 数据编辑窗口是SPSS的主程序窗口，可以同时打开多个数据编辑窗口。
- ✓ 数据编辑窗口的主要功能是：
定义SPSS数据的结构、录入编辑和管理待分析的数据。
- ✓ 这些数据通常以SPSS数据文件的形式保存，文件扩展名为.sav。

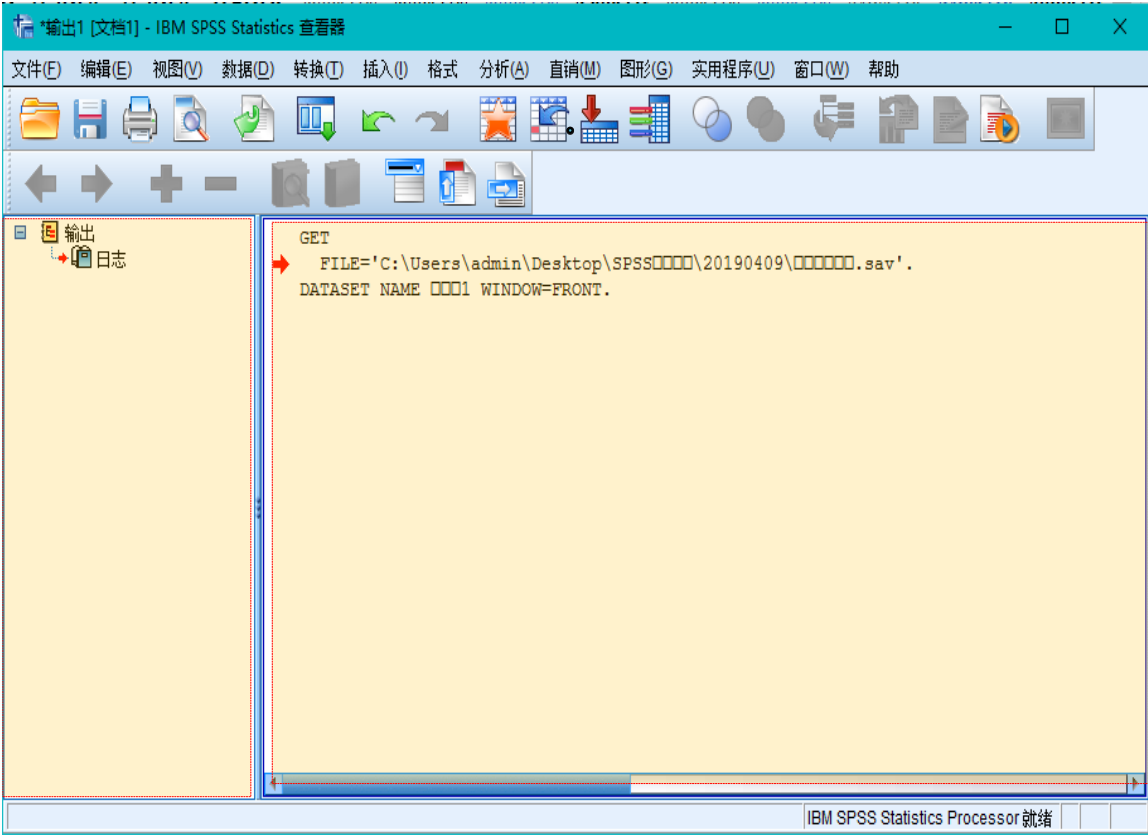
结果输出窗口

窗口主菜单

与编辑窗口主菜单选项大致相同
有两项独有的菜单，分别是插入和格式

分析结果显示区

目录区：分析结果的目录
内容区：分析结果的详细报告



工具栏

保留数据编辑窗口中的某些功能按钮
还增添了一些自己特有的功能按钮

系统状态显示区

显示系统的当前运行状态

- ✓ 结果输出窗口是SPSS的主要窗口，允许同时创建或打开多个输出窗口。
- ✓ 结果窗口的主要功能是：
显示、管理SPSS统计分析结果、报表及图形。
- ✓ 这些数据通常以SPSS输出文件的形式保存，文件扩展名为.spv。

SPSS基本运行方式

SPSS为用户提供了以下三种基本运行方式

分别适合于不同的用户和不同的统计分析要求

今天的培训内容都基于完全窗口菜单方式



完全窗口菜单方式

- 所有分析操作过程都通过菜单和按钮及对话框方式进行
- 适用于一般分析和SPSS初学者



程序运行方式

- 手工编写SPSS命令程序，一次性提交计算机运行
- 适用于大规模的分析工作和熟练的SPSS程序员



菜单程序混合运行方式

- 使用菜单的时候同时编辑SPSS程序
- 适用于熟练的SPSS程序员

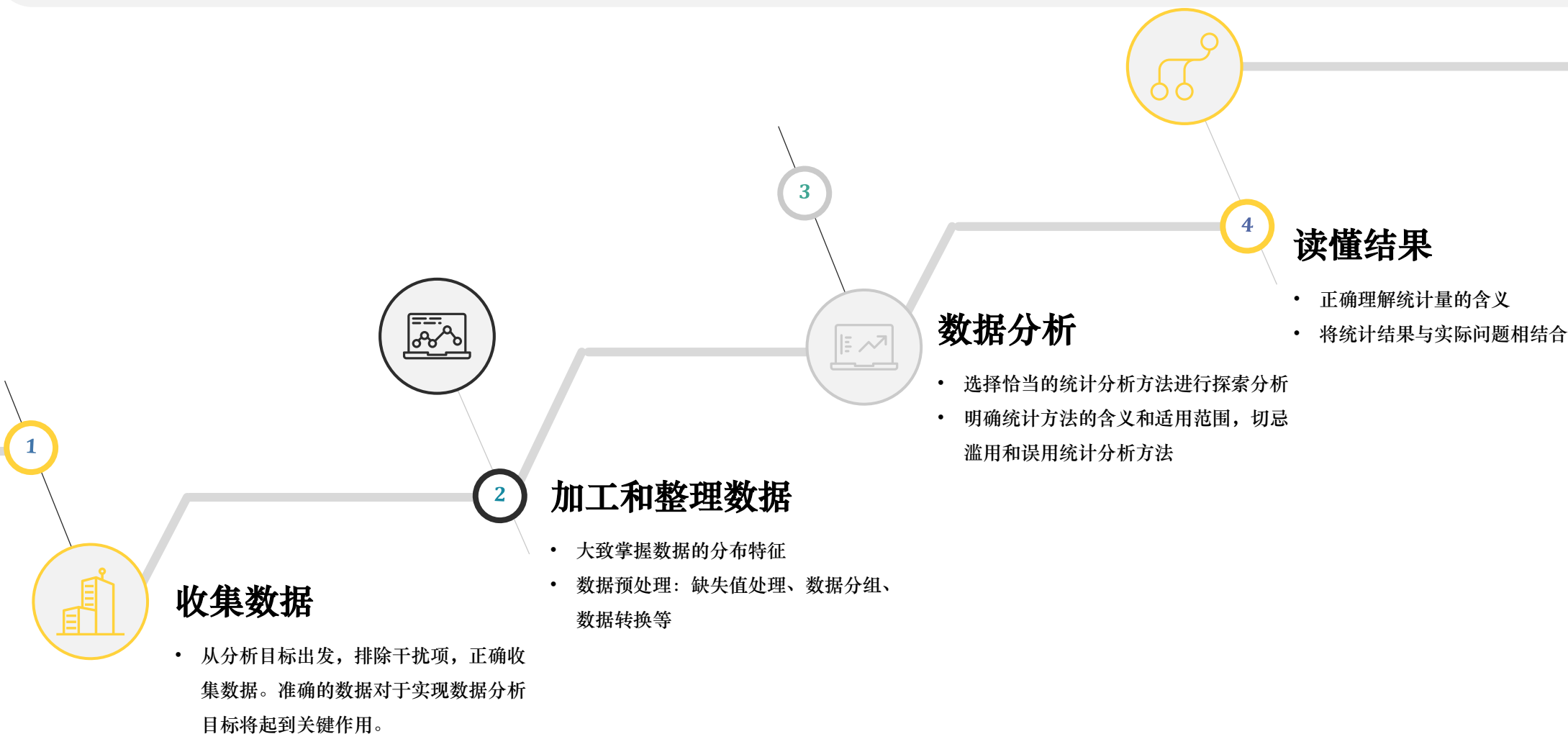


/02

数据分析的一般步骤

数据分析的一般步骤

明确数据分析目标是数据分析的出发点！



利用SPSS进行数据分析的一般步骤





/03

问卷数据分析案例

以2016年中国流动人口动态监测调查数据为例

数据样例：

<http://hdl.handle.net/20.500.12291/10215>

案例介绍

中国流动人口动态监测调查数据

国家卫生健康委自2009年起一年一度大规模全国性流动人口抽样调查数据，覆盖全国31个省（区、市）和新疆生产建设兵团中流动人口较为集中的流入地，内容涉及流动人口基本信息、流动范围和趋向、就业和社会保障、收支和居住、基本公共卫生服务、婚育和计划生育服务管理、子女流动和教育等。



1

- 数据结构
- 问卷录入（多选题、矩阵题）

2

- 数据预处理（变量计算、数据选取、分类汇总、数据分组）
- 基本统计分析（频数分析、基本描述统计量、多选项分析）

3

- t 检验
- 方差分析

数据结构

对每列变量及其相关属性的描述



变量名（Name）
变量存取的唯一标志



变量名标签（Label）
对变量名的一些解释说明，增强分析结果的可视性和分析结果的可读性



变量值标签（Values）
对变量取值含义的一些解释说明，对定类型和定距型数据尤为重要

*2016_a.sav [数据集1] - IBM SPSS Statistics 数据编辑器

	名称	类型	宽度	小数	标签	值	缺失	列	对齐
1	ID	数值(N)	12	0	顺序编码	无	无	8	右
2	C1	字符串	16	0	现居住地所在省	无	无	16	左
3	C2	字符串	20	0	现居住地所在(地区)	无	无	20	左
4	C3	字符串	20	0	现居住地所在县(市、区)	无	无	20	左
5	C6	数值(N)	12	0	样本点编码	无	无	8	右
6	C7	数值(N)	12	0	样本点类型	{1, 居委会}...	无	8	右
7	C8	数值(N)	12	0	被访者编码	无	无	8	右
8	Q100	数值(N)	12	0	同住的家庭成员人数	无	无	8	右
9	Q101D1	数值(N)	12	0	成员序列号	无	无	8	右
10	Q101A1	数值(N)	12	0	与被访者关系	{1, 本人}...	无	8	右
11	Q101B1	数值(N)	12	0	性别	{1, 男}...	无	8	右
12	Q101C1Y	数值(N)	12	0	出生年	无	无	8	右
13	Q101C1M	数值(N)	12	0	出生月	无	无	8	右
14	Q101D1	数值(N)	12	0	民族	{1, 汉族}...	无	8	右
15	Q101E1	数值(N)	12	0	受教育程度	{1, 未上过学}...	无	8	右
16	Q101F1	数值(N)	12	0	户口性质	{1, 农业}...	无	8	右
17	Q101G1	数值(N)	12	0	婚姻状况	{1, 未婚}...	无	8	右
18	Q101H1	数值(N)	12	0	是否中共党员	{1, 是}...	无	8	右
19	Q101I1	数值(N)	12	0	是否本地户籍人口	{1, 是}...	无	8	右
20	Q101J1	数值(N)	12	0	户籍地省份	无	无	8	右
21	Q101K1	数值(N)	12	0	现居住地	{1, 本地}...	无	8	右
22	Q101L1	数值(N)	12	0	本次流动范围	{1, 跨省}...	无	8	右
23	Q101M1Y	数值(N)	12	0	本次流动年份	无	无	8	右
24	Q101M1M	数值(N)	12	0	本次流动月份	无	无	8	右
25	Q101N1	数值(N)	12	0	本次流动原因	{1, 务工/工}...	无	8	右
26	Q101D2	数值(N)	12	0	成员序列号	无	无	8	右
27	Q101A2	数值(N)	12	0	与被访者关系	{1, 本人}...	无	8	右
28	Q101B2	数值(N)	12	0	性别	{1, 男}...	无	8	右
29	Q101C2Y	数值(N)	12	0	出生年	无	无	8	右
30	Q101C2M	数值(N)	12	0	出生月	无	无	8	右
31	Q101D2	数值(N)	12	0	民族	{1, 汉族}...	无	8	右
32	Q101E2	数值(N)	12	0	受教育程度	{1, 未上过学}...	无	8	右
33	Q101F2	数值(N)	12	0	户口性质	{1, 农业}...	无	8	右

数据视图 变量视图

IBM SPSS Statistics Processor 就绪

数据结构

对每列变量及其相关属性的描述



数据类型（Type）、列宽（Width）、小数（Decimals）

三种基本数据类型：

数值型（Numeric）、字符型（String）、日期型（Date）

相应的类型有默认的列宽或小数位宽



缺失值（Missing）

漏填数据、明显错误或不合理的数据等

用户缺失值：指定某个特定值为缺失值

系统缺失值：数值型：点（.）；字符型：空



度量标准（Measure）

Scale：定距型数据；Ordinal：定序型数据；Nominal：定类型数据（无顺序）

	名称	类型	宽度	小数	标签	值	缺失	列	对齐
1	ID	数值(N)	12	0	顺序编码	无	无	8	右
2	C1	字符串	16	0	现居住地所在省	无	无	16	左
3	C2	字符串	20	0	现居住地所在(地区)	无	无	20	左
4	C3	字符串	20	0	现居住地所在(市、区)	无	无	20	左
5	C6	数值(N)	12	0	样本点编码	无	无	8	右
6	C7	数值(N)	12	0	样本点类型	{1, 居委会}...	无	8	右
7	C8	数值(N)	12	0	被访者编码	无	无	8	右
8	Q100	数值(N)	12	0	同住的家庭成员人数	无	无	8	右
9	Q101D1	数值(N)	12	0	成员序列号	无	无	8	右
10	Q101A1	数值(N)	12	0	与被访者关系	{1, 本人}...	无	8	右
11	Q101B1	数值(N)	12	0	性别	{1, 男}...	无	8	右
12	Q101C1Y	数值(N)	12	0	出生年	无	无	8	右
13	Q101C1M	数值(N)	12	0	出生月	无	无	8	右
14	Q101D1	数值(N)	12	0	民族	{1, 汉族}...	无	8	右
15	Q101E1	数值(N)	12	0	受教育程度	{1, 未上过学}...	无	8	右
16	Q101F1	数值(N)	12	0	户口性质	{1, 农业}...	无	8	右
17	Q101G1	数值(N)	12	0	婚姻状况	{1, 未婚}...	无	8	右
18	Q101H1	数值(N)	12	0	是否中共党员	{1, 是}...	无	8	右
19	Q101I1	数值(N)	12	0	是否本地户籍人口	{1, 是}...	无	8	右
20	Q101J1	数值(N)	12	0	户籍地省份	无	无	8	右
21	Q101K1	数值(N)	12	0	现居住地	{1, 本地}...	无	8	右
22	Q101L1	数值(N)	12	0	本次流动范围	{1, 跨省}...	无	8	右
23	Q101M1Y	数值(N)	12	0	本次流动年份	无	无	8	右
24	Q101M1M	数值(N)	12	0	本次流动月份	无	无	8	右
25	Q101N1	数值(N)	12	0	本次流动原因	{1, 务工/工}...	无	8	右
26	Q101D2	数值(N)	12	0	成员序列号	无	无	8	右
27	Q101A2	数值(N)	12	0	与被访者关系	{1, 本人}...	无	8	右
28	Q101B2	数值(N)	12	0	性别	{1, 男}...	无	8	右
29	Q101C2Y	数值(N)	12	0	出生年	无	无	8	右
30	Q101C2M	数值(N)	12	0	出生月	无	无	8	右
31	Q101D2	数值(N)	12	0	民族	{1, 汉族}...	无	8	右
32	Q101E2	数值(N)	12	0	受教育程度	{1, 未上过学}...	无	8	右
33	Q101F2	数值(N)	12	0	户口性质	{1, 农业}...	无	8	右

多选题

二分法

将多选题的每个答案设为一个变量
每个变量具有两个可能值：
是/否、有/无、选中/未选中

多重分类法

用多个变量来记录问题的答案
变量值标签一致，均包含所有选项
通常用于“最多选几项”或者“多选答案存在顺序”时使用

名称	类型	宽度	小数	标签	值
Q303A	数值(N)	12	0	您是否在本地产购买住房	{1, 是}...
Q303B	数值(N)	12	0	您是否在本地产-区政府所在地购买住房	{1, 是}...
Q303C	数值(N)	12	0	您是否在本地产-乡政府所在地购买住房	{1, 是}...
Q303D	数值(N)	12	0	您是否在本地产-村购买住房	{1, 是}...
Q303E	数值(N)	12	0	您是否在其他地方购买住房	{1, 是}...

多选题

值标签(V) X

值标签(V)

值(U):

拼写(S)...

标签(L):

添加(A)

更改(C)

删除(R)

1 = "笔记本电脑"
2 = "台式电脑"
3 = "数码相机"
4 = "数码摄像机"
5 = "手机"
6 = "MP3/MP4"
7 = "电子书"
8 = "其他"

确定

取消

帮助

多重分类法

用多个变量来记录问题的答案
变量值标签一致，均包含所有选项
通常用于“最多选几项”或者“多选答案存在顺序”时使用

名称	类型	宽度	小数	标签	值
B2_1	数值(N)	8	0	请问您未来打算优先购买的数码产品有哪些？（最多三项）	{1, 笔记本电脑}...
B2_2	数值(N)	8	0	请问您未来打算优先购买的数码产品有哪些？（最多三项）	{1, 笔记本电脑}...
B2_3	数值(N)	8	0	请问您未来打算优先购买的数码产品有哪些？（最多三项）	{1, 笔记本电脑}...

矩阵题

219 您目前参加下列何种社会保障？

社会保障	1. 是否参保	2. 在何处参保
	1 是	1 本地
	2 否（跳问下一行）	2 户籍地
	3 不清楚（跳问下一行）	3 其他地方
A 养老保险（含新农保、养老金等）	<input type="checkbox"/>	<input type="checkbox"/>
B 失业保险	<input type="checkbox"/>	<input type="checkbox"/>
C 工伤保险	<input type="checkbox"/>	<input type="checkbox"/>
D 生育保险	<input type="checkbox"/>	<input type="checkbox"/>
E 住房公积金	<input type="checkbox"/>	<input type="checkbox"/>

名称	类型	宽度	小数	标签	值
Q219A1	数值(N)	12	0	您是否参加养老保险(含新农保、养老金等)	{1, 是}...
Q219A2	数值(N)	12	0	您在何地参加养老保险(含新农保、养老金等)	{1, 本地}...
Q219B1	数值(N)	12	0	您是否参加失业保险	{1, 是}...
Q219B2	数值(N)	12	0	您在何地参加失业保险	{1, 本地}...
Q219C1	数值(N)	12	0	您是否参加工伤保险	{1, 是}...
Q219C2	数值(N)	12	0	您在何地参加工伤保险	{1, 本地}...
Q219D1	数值(N)	12	0	您是否参加生育保险	{1, 是}...
Q219D2	数值(N)	12	0	您在何地参加生育保险	{1, 本地}...

数据预处理



变量计算

- SPSS算术表达式
- SPSS条件表达式
- SPSS函数



数据选取

- 按指定条件选取
- 随机选取



分类汇总

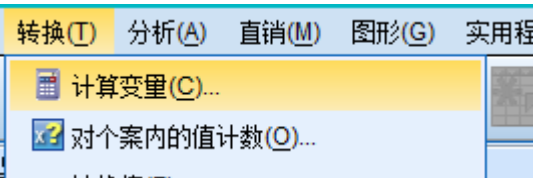
- 按照某分类变量进行分类计算
- 根据多个分类变量汇总计算时称为多重分类汇总



数据分组

- 根据需要将数据按照某种标准重新划分为不同的组别

变量计算



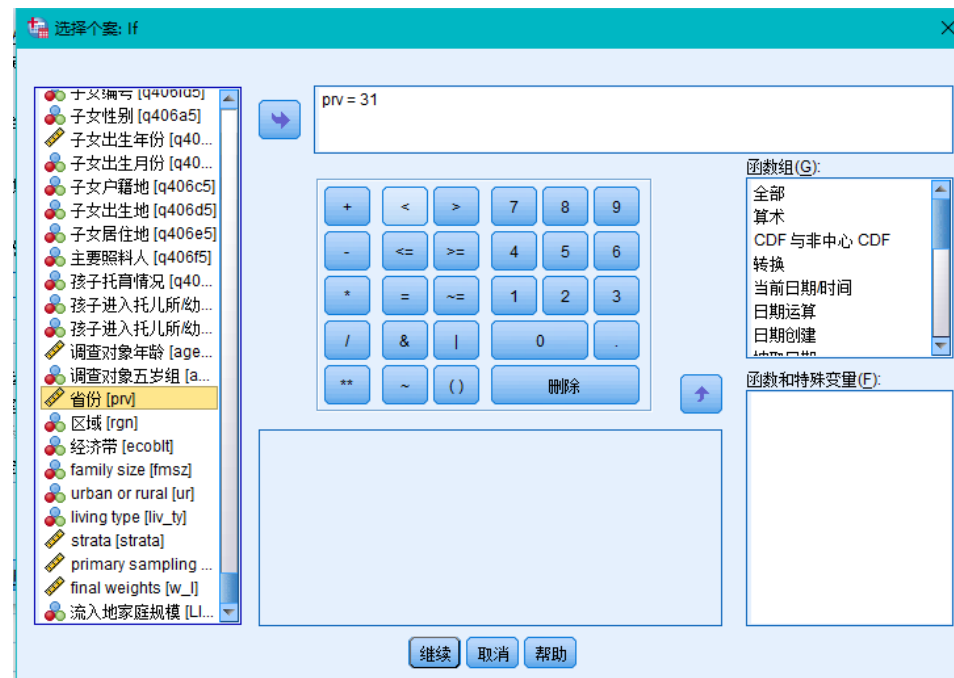
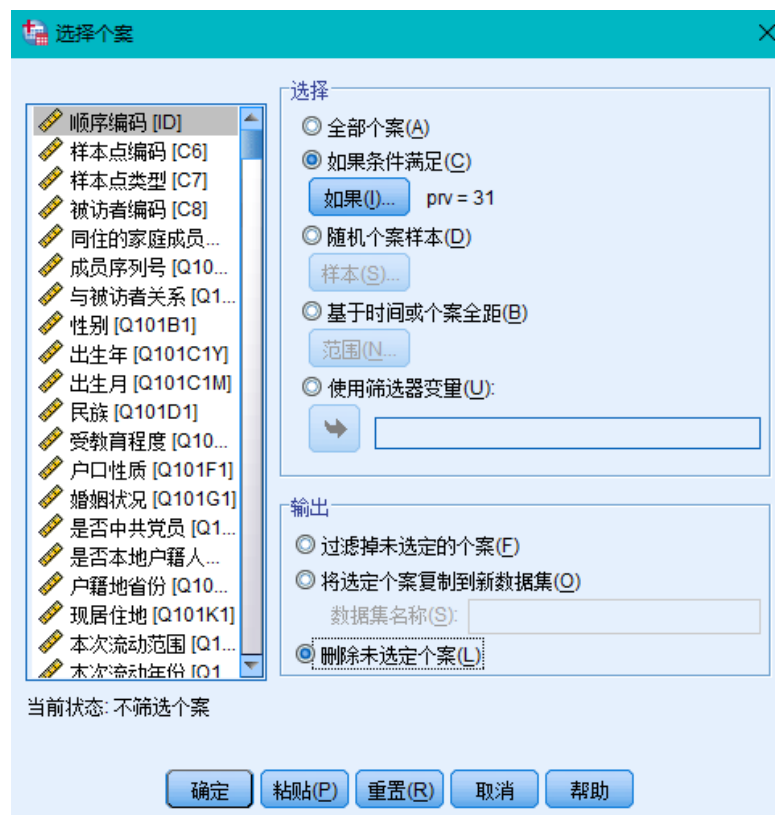
根据出生年、月计算年龄



数据选取



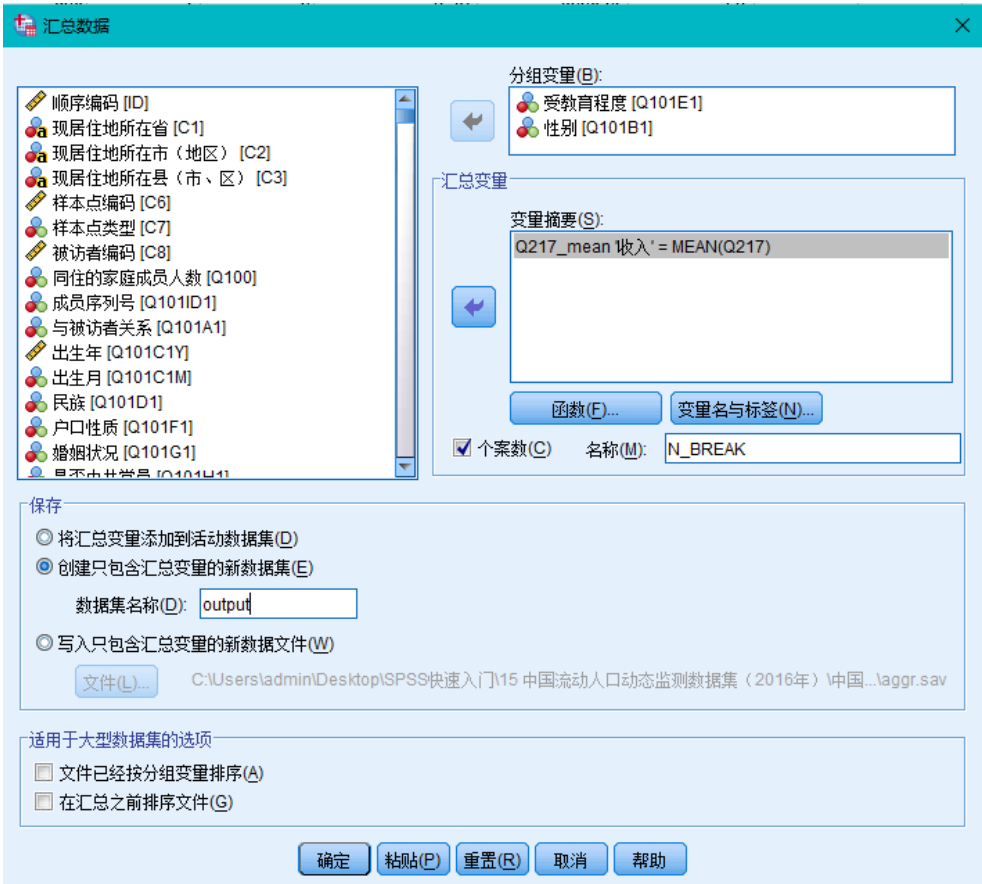
选取上海地区的样本



分类汇总

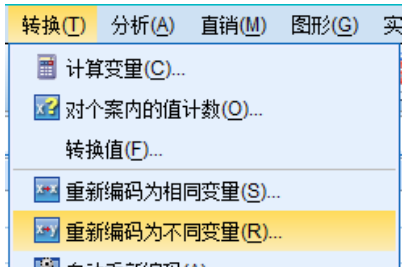


按受教育程度进行分类，计算流动人口的上月平均收入



Q101E1	平均收入	样本量
1	3324.58	117
2	3765.08	686
3	4776.13	2724
4	5423.06	1434
5	7906.62	875
6	9567.44	1019
7	11826.76	145

数据分组



年龄分组

重新编码为其他变量

数字变量 -> 输出变量(V):
age -> agex

输出变量
名称(N):
agex
标签(L):
年龄分组
更改(H)

旧值和新值(O)...
如果(I)... (可选的个案选择条件)

确定 粘贴(P) 重置(R) 取消 帮助

重新编码到其他变量: 旧值和新值

旧值
☒ 值(V):
☐ 系统缺失(S)
☐ 系统或用户缺失(U)
☐ 范围:
到(T)
☐ 范围, 从最低到值(G):
☐ 范围, 从值到最高(E):
☐ 所有其他值(O)

新值
☒ 值(L):
☐ 系统缺失(Y)
☐ 复制旧值(P)

旧 -> 新(O):
Lowest thru 19 -> 1
20 thru 39 -> 2
40 thru 59 -> 3
60 thru Highest -> 4

添加(A) 更改(C) 删除(R)

☐ 输出变量为字符串(B)
☐ 将数值字符串移动为数值(U)

继续 取消 帮助

	Q208X_mean	birthday	today
1	41.94	Dec 1986	Jun 2020
1	40.78	Jul 1984	Jun 2020
3	42.97	Oct 1982	Jun 2020
3	49.75	Sep 1982	Jun 2020
3	51.75	Aug 1973	Jun 2020
3	51.08	Oct 1971	Jun 2020
3	50.75	Aug 1971	Jun 2020
3	51.08	Jan 1972	Jun 2020
2	48.59	Nov 1997	Jun 2020

2	4	1	3	16
2	3	1	3	16
2	3	1	3	16
2	3	1	3	16
2	4	1	3	16
2	1	1	3	16
2	2	1	3	16
2	2	1	3	16
2	2	1	3	16
2	2	1	3	16
2	3	1	3	16
2	3	1	3	16

基本统计分析

掌握数据的基本统计特征
把握数据的总体分布形态



频数分析

- 编制频数分布表
- 绘制统计图
- 交叉列联表



计算描述统计量

- 集中趋势：均值、中位数、众数
- 离散程度：样本标准差、样本方差
- 分布形态：偏度系数、峰度系数



多选项分析

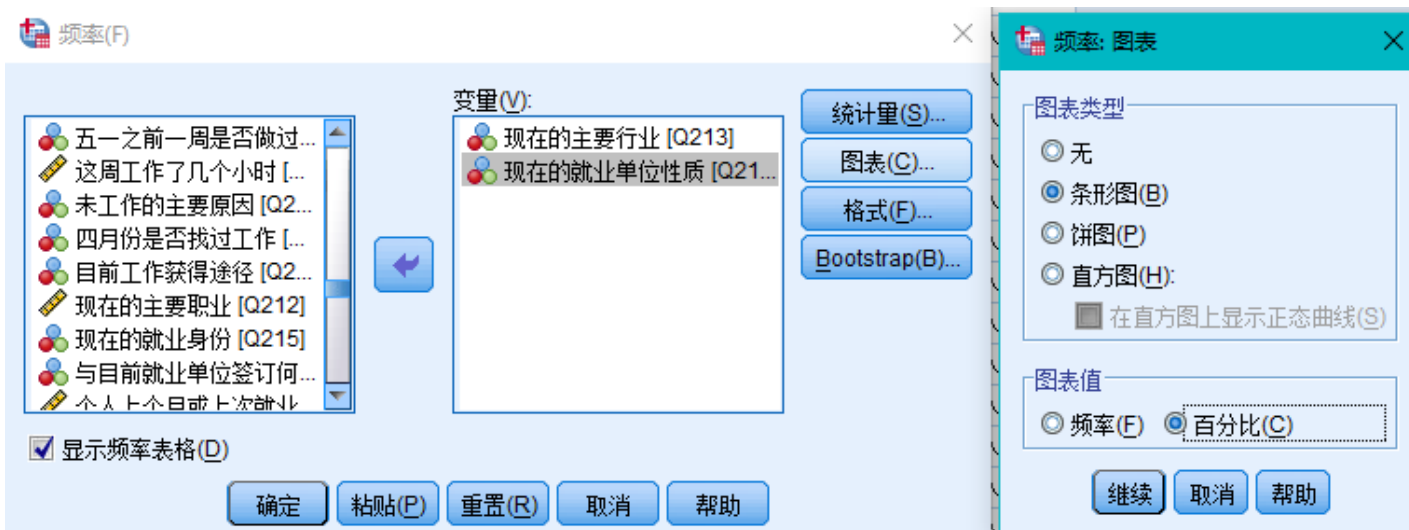
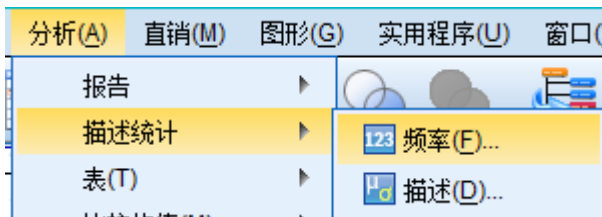
- 按二分法或多重分类法设置变量
- 多选项频数分析
- 多选项交叉分组下的频数分析



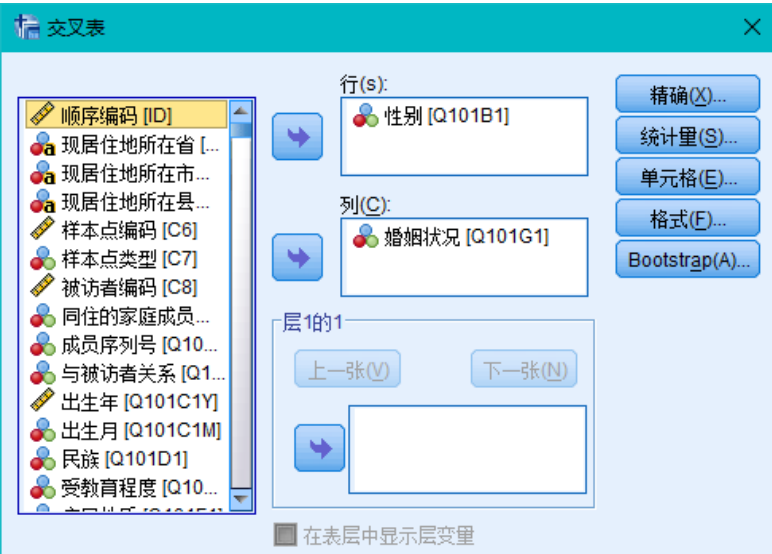
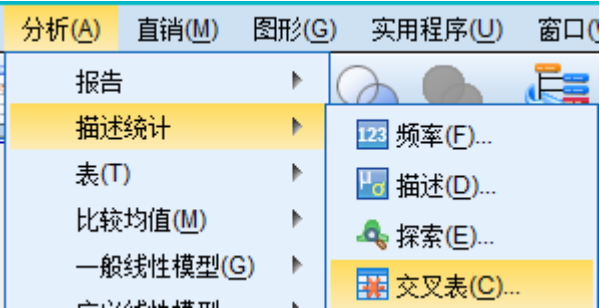
其他探索性分析

频数分析

查看就业单位性质的频数分布



交叉列联表



性别、婚姻状况交叉列联表

卡方检验

	值	df	渐进 Sig. (双侧)
Pearson 卡方	28.286 ^a	5	.000
似然比	29.559	5	.000
线性和线性组合	12.218	1	.000
有效案例中的 N	7000		

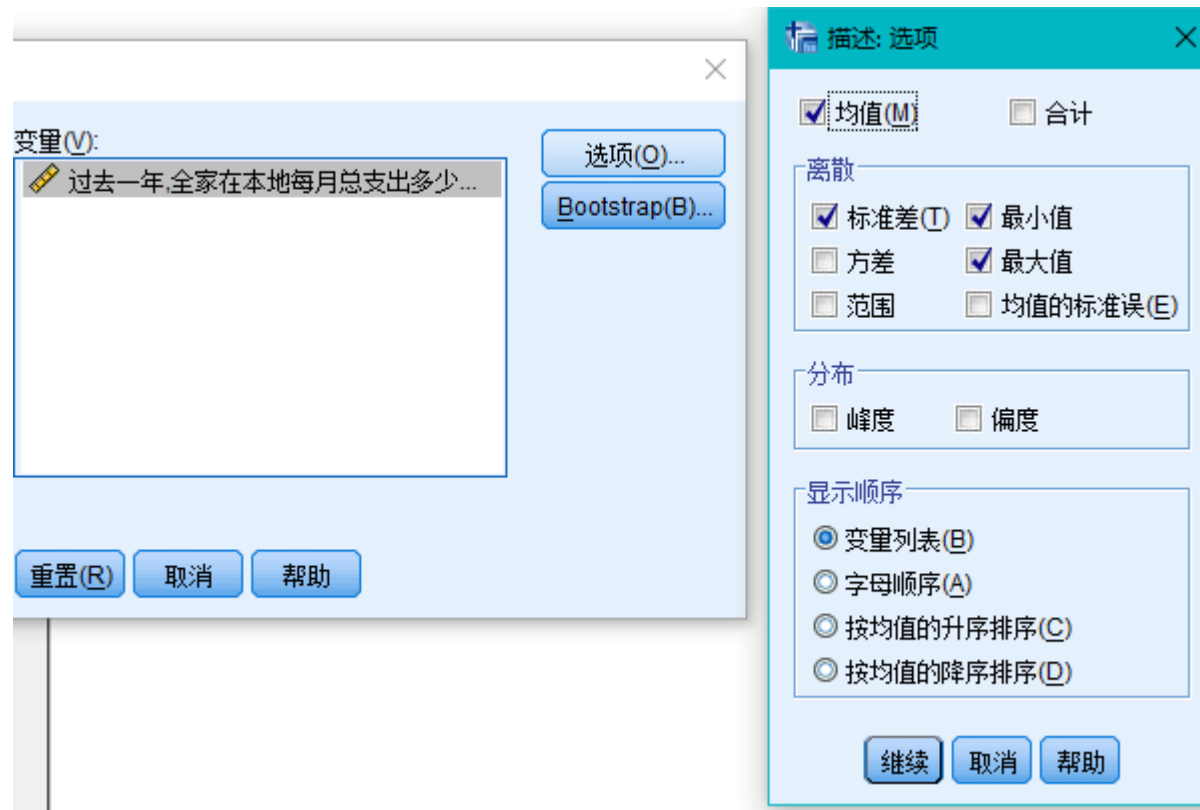
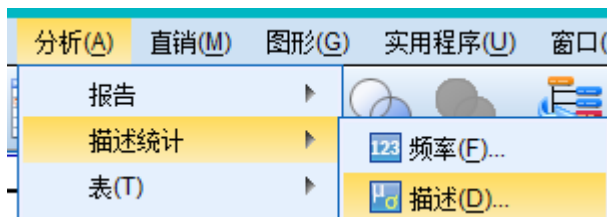
a. 0 单元格(0.0%) 的期望计数少于 5。最小期望计数为 17.77。

性别* 婚姻状况 交叉制表

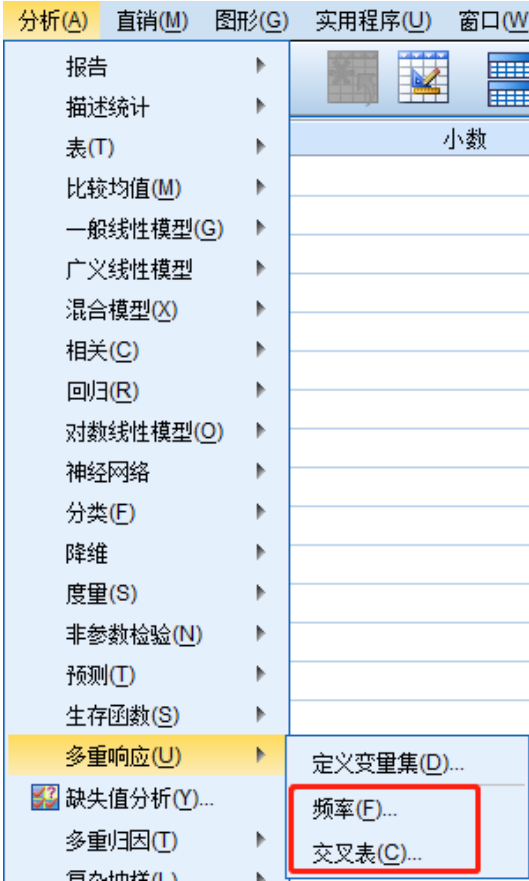
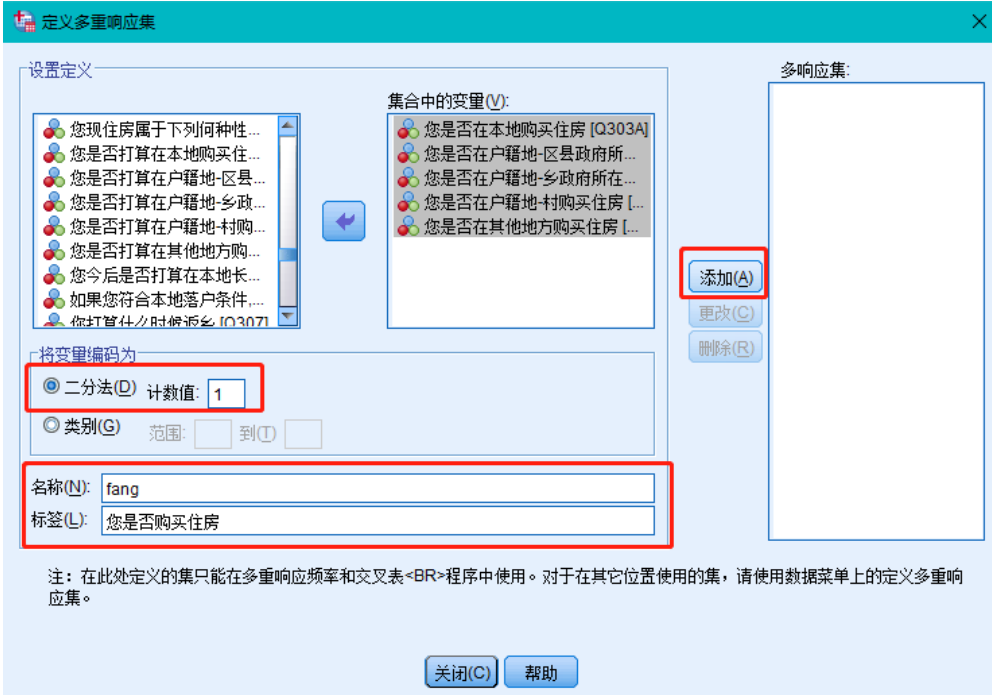
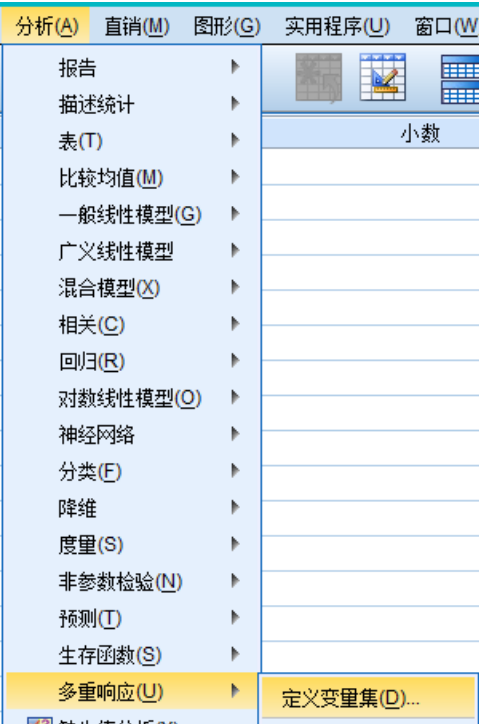
计数

		婚姻状况						合计
		未婚	初婚	再婚	离婚	丧偶	同居	
性别	男	365	2869	51	37	7	32	3361
	女	405	3005	89	57	30	53	3639
合计		770	5874	140	94	37	85	7000

计算描述统计量



多选项分析



推断统计与参数检验

推断统计

- 根据样本数据推断总体特征
总体数据无法全部收集
采集数据需要大量投入



参数检验

- 参数检验是推断统计的重要组成部分
- 当总体分布已知（如总体为正态分布），根据样本数据对总体分布的统计参数进行推断
- 也可对两个或多个总体的参数进行比较

t 检验

假设检验的一般步骤：

- 第一步：提出原假设
- 第二步：选择检验统计量
- 第三步：计算检验统计量观测值发生的概率
- 第四步：给定显著性水平 α ，并作出统计决策

利用SPSS进行假设检验时，首先需明确原假设，第二步和第三步是SPSS自动完成的，最终决策需要人工判定。



单样本t检验

- 对总体均值的假设检验
- 前提：样本来自的总体应服从或近似服从正态分布
- 检验统计量为t统计量



两独立样本t检验

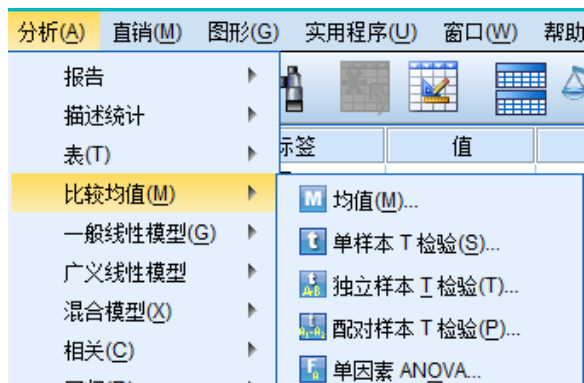
- 推断两个总体的均值是否存在显著差异
- 前提1：样本来自的总体应服从或近似服从正态分布
- 前提2：两样本相互独立，样本量可以不等



两配对样本t检验

- 配对样本通常具有两个特征：
 - 一、两组样本量相同
 - 二、两组样本观测值先后顺序一一对应

t 检验



组统计量

	性别	N	均值	标准差	均值的标准误
总共流动次数	男	88087	1.40	1.193	.004
	女	80911	1.27	.811	.003

男性和女性流动人口的流动次数存在显著差异

独立样本检验

		方差方程的 Levene 检验		均值方程的 t 检验					
		F	Sig.	t	df	Sig.(双侧)	均值差值	标准误差值	差分的 95% 置信区间
总共流动次数	假设方差相等	1862.863	.000	25.833	168996	.000	.129	.005	.120 .139
	假设方差不相等			26.240	156011.381	.000	.129	.005	.120 .139

方差分析

- 通过推断控制变量各水平下各观测变量的总体均值是否存在显著差异，分析控制变量是否给观测变量带来了显著影响。
- 基本假设前提：
 - 一、观测变量各总体服从正态分布
 - 二、观测变量各总体方差相同



单因素方差分析

- 仅研究单个因素对观测变量的影响
- 检验统计量为F统计量
- 方差齐性检验
- 多重比较检验



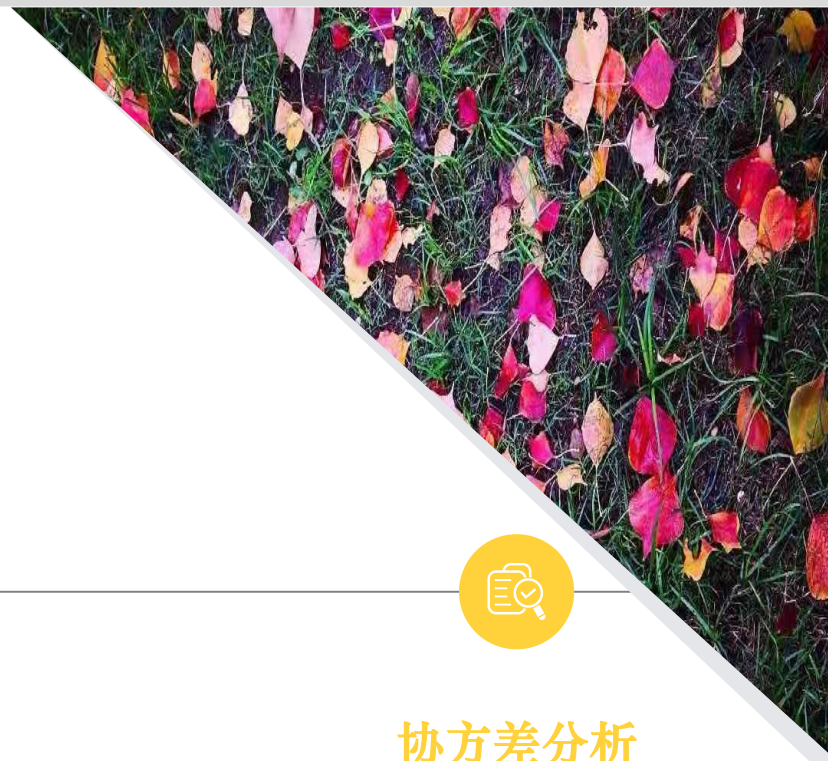
多因素方差分析

- 研究多个因素对观测变量的影响
- 控制变量独立作用的影响
- 控制变量交互作用的影响
- 随机因素的影响

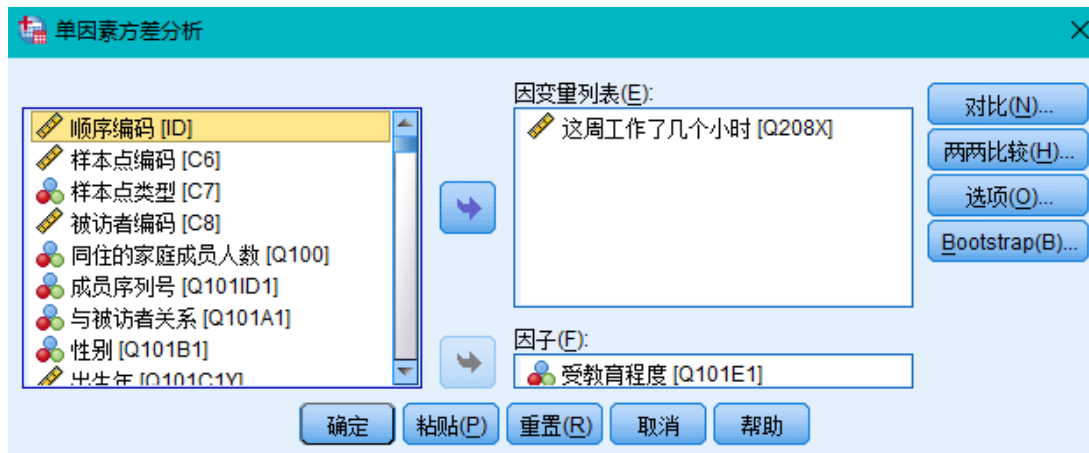
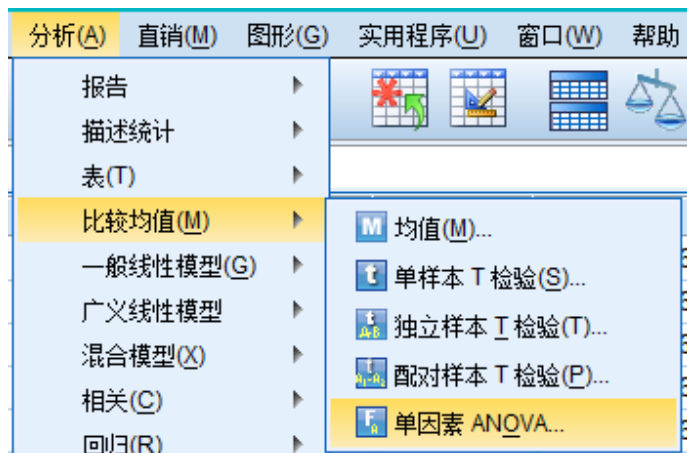


协方差分析

- 协变量：人为因素很难控制的因素



方差分析



方差分析

方差齐性检验

这周工作了几个小时

Levene 统计量	df1	df2	显著性
189.919	6	5899	.000

单因素方差分析

这周工作了几个小时

	平方和	df	均方	F	显著性
组间	88044.778	6	14674.130	88.946	.000
组内	973205.502	5899	164.978		
总数	1061250.280	5905			

受教育程度对工作时长是否产生显著影响

多重比较

因变量: 这周工作了几个小时

(I) 受教育程度	(J) 受教育程度	均值差 (I-J)	标准误	显著性	95% 置信区间	
					下限	上限
LSD 未上过学	小学	1.032	1.609	.521	-2.12	4.18
	初中	2.126	1.537	.167	-.89	5.14
	高中/中专	5.760*	1.559	.000	2.70	8.82
	大学专科	9.970*	1.585	.000	6.86	13.08
	大学本科	11.202*	1.572	.000	8.12	14.28
	研究生	10.927*	1.860	.000	7.28	14.57
小学	未上过学	-1.032	1.609	.521	-4.18	2.12
	初中	1.094	.606	.071	-.09	2.28
	高中/中专	4.729*	.660	.000	3.43	6.02
	大学专科	8.938*	.720	.000	7.53	10.35
	大学本科	10.170*	.691	.000	8.82	11.52
	研究生	9.895*	1.211	.000	7.52	12.27
初中	未上过学	-2.126	1.537	.167	-5.14	.89
	小学	-1.094	.606	.071	-2.28	.09
	高中/中专	3.635*	.460	.000	2.73	4.54
	大学专科	7.844*	.542	.000	6.78	8.91
	大学本科	9.076*	.503	.000	8.09	10.06
	研究生	8.801*	1.114	.000	6.62	10.99



Thanks



復旦大學 大 数 据 研 究 院
人文社会科学数据研究所
Institute for Humanities and Social Science Data