

```

1      *=====
2      *第2讲 数据文件建立和管理（初级篇）
3      *=====
4
5      *-----
6      *->2.1数据文件的建立和读取
7      *-----
8      *-2.1.1.直接录入（适用于小样本少变量的数据文件新建）
9      *    直接在stata中录入：打开程序，调用数据编辑窗口，直接录入数据，如excel中操作。
10     *    调用数据窗口方式：
11     *        (a)在command窗口中输入edit命令
12     *        (b)点工具栏上的编辑文件图标
13     *    命令行输入：在command窗口，用input命令。
14     *操作实例->
15     input  x  y
16     1    2
17     2    3
18     end
19     *-2.1.2.读取dta数据文件
20     *    鼠标或菜单操作：只能读取本地磁盘上的Stata数据文件
21     *    直接双击stata数据文件
22     *    菜单操作：在工具栏上直接点击打开文件图标，打开指定文件，或File>Open
23     *    use命令
24     *    语法：use [varlist] [if] [in] using filename [,clear nolabel]
25     *操作实例->
26     cd "C:\data"
27     *    读取数据文件的指定变量信息
28     use make mpg using auto.dta, clear // clear: 在读取新数据时清除程序已读数据
29     *    读取数据文件的全部变量信息
30     use auto.dta, clear
31     use auto.dta, clear nolabel // nolabel: 在读取数据时清除数据中变量的值标签
32     *    读取数据文件中符合指定条件的样本信息
33     use auto.dta if price>7000, clear
34     *逻辑符号
35     *    ==: 等于, !=: 不等于
36     *    > : 大于, >=: 大于等于
37     *    < : 小于, <=: 小于等于
38     *    & : 且 , | : 或
39     *    读取数据文件中符合指定范围的样本信息
40     use auto.dta in 1/20, clear
41     * 范围语句的定义:
42     *    in #: 第#个观测值
43     *    in #1/#2: 从#1个观测值到第#2个观测值
44     *    in f/#: 从第1个观测值到第#个观测值
45     *    in #/l: 从第#个观测值到最后一个观测值
46     *    use命令的扩展: sysuse/webuse/netuse
47     *-2.1.3.导入其它格式的数据文件
48     *    支持导入的数据类型包括: (1) Excel数据: *.xls, *.xlsx;
49     // (2) 用spreadsheet建立的ASCII(txt)数据: *.raw, *.txt, *.csv;
50     // (3) 固定列宽的ASCII(txt)数据: *.dct;
51     // (4) 以dictionary格式建的ASCII(txt)数据: *.dct, *.raw;
52     // (5) 无固定格式的ASCII(txt)数据: *.txt, *.raw;
53     // (6) SAS XPORT数据: *.xpt;
54     // (7) ODBC数据源: 包括access数据源, *.mdb, dBase数据源, *.dbf;
55     // (8) xml数据: *.xml
56     *    import excel命令: 导入Excel数据, 文件后缀名是*.xls, *.xlsx
57     *    语法1: import excel [using] filename [, import_excel_options] //导入整个Excel数据工作表
58     *    语法2: import excel extvarlist using filename [, import_excel_options]
59     //导入Excel数据工作表的部分变量
60     *操作实例->
61     cd "C:\data"
62     sysuse auto
63     export excel auto, firstrow(variables)
64     import excel auto.xls, firstrow clear
65     describe
66     import excel auto.xls, cellrange(:D70) firstrow clear
67     describe
68     import excel make mpg weight price using auto.xls, firstrow clear

```

```

68 describe
69 import excel make=A mpg=B price=D using auto.xls, firstrow clear
70 describe
71 * insheet命令: 导入以tab-/逗号 (comma-) /指定字符分隔的数据, 文件后缀名通常是*.raw, *.txt,
  *.csv。
72 //空格分隔的数据无法直接用insheet命令, 需要在参数设定delimiter(" ")。
73 * 语法: insheet [varlist] using filename [, options]
74 *操作实例->
75 insheet using auto.raw, tab
76 insheet using auto.txt, comma clear
77 insheet mpg price using auto.raw, comma
78 insheet mpg price using auto.raw, delimiter("/")
79 * infile命令: 导入以空格 (space) -/tab-/逗号 (comma-) 分隔的数据, 文件后缀名通常.dct, *.raw,
  *.txt。注意: 在raw, txt格式的数据文件导入, 一定要在命令中写出导入后数据的对应变量名。
80 *操作实例->
81 infile v1 v2 v3 using d21.txt, clear
82 infile x1 x2 using infile.txt, clear
83
84 *-----
85 *->2.2数据的保存或导出
86 *-----
87 *-2.2.1.数据保存
88 * save命令
89 * 语法: save [filename] [, replace/nolabel/orphans /emptyok]
90 *操作实例->
91 save auto1a, replace
92 save auto1b, replace nolabel
93 save auto1c, replace orphans //保存所有值标签
94 *-2.2.2.数据导出
95 * export excel/outfile/outsheet/fdasave/xmlsave命令
96
97 *-----
98 *->2.3数据清理的常见操作
99 *-----
100 *-2.3.1.变量属性, 如名称、标签、值标签
101 * 变量名: 由英文字符、数字、中文字符和_组成, 最多不超过32个字母。
102 * 字母大小写表示的含义不同!!!
103 * 数字不可以单独做变量名, 也不可放在变量名的最前面。
104 * 建议不要使用_作为变量的第一个字母。因为许多stata的内部变量都是以_开头的, 如, _n, _N,
  _cons, _b等等。
105 * 标签: 对变量的含义进行解释
106 * 值标签: 对分类变量值的含义进行解释
107 *例子: 性别: 女=0, 男=1
108 *-2.3.2.变量属性的显示
109 * describe命令
110 * 语法1: describe [varlist], [simple/short]
111 * 语法2: describe [varlist] using filename, [simple/short]
112 *操作实例->
113 sysuse auto, clear
114 describe
115 describe make price mpg
116 describe course using "c:\data\myscore.dta"
117 *-2.3.3.变量属性的修改
118 * rename命令: 变量名更改
119 * 语法1: rename old_var new_var
120 * 语法2: rename (old1 old2 ...) (new1 new2 ...) [, options1]
121 *操作实例->
122 sysuse auto, clear
123 rename price pri
124 rename (make price mpg) (mk pr mp)
125 rename t* t1*
126 * label命令: 定义和管理数据、变量或变量值的标签
127 * 语法1: label data "label" //定义数据标签
128 * 语法2: label var varname "label" //定义变量标签
129 * 语法3:
130 * label define lblname # "label" [# "label" ...] [, add modify replace nofix] //定义值标签
131 * label values varlist [lblname|.] [, nofix] //附加值标签到指定分类变量上
132 *操作实例->

```

```

133 sysuse auto, clear
134 label data "auto in American" //数据标签
135 label var foreign "car type" //变量标签
136 label define origin 0 "domestic" 1 "foreign" //值标签
137 label values foreign origin //附加值标签origin到分类变量foreign上
138 * 语法4: label dir //显示值标签名
139 * 语法5: label list //显示值标签名和内容
140 * 语法6: label drop //剔除值标签
141 * 语法7: label copy //复制值标签
142 * 语法8: label save //保存值标签到do文件里
143 *-2.3.4.变量的存储类型
144 * 整数（数值型变量）的存储类型：
145 * byte 字节型 (-100, +100)
146 * int 一般整数型 (-32000, +32000)
147 * long 长整数型 (-2.14*10^10, +2.14*10^10)，即，正负21亿
148 * 小数（数值型变量）的存储类型：
149 * float 浮点型 8位有效数字
150 * double 双精度 16位有效数字
151 * 字符型变量的存储类型: str#
152 *变量信息由字母、特殊符号和数字组成
153 *保存为str格式，str后面的数字代表最大字符长度。比如，str18
154 *用在英文状态下的引号""标注
155 * 日期型变量，如19870815或15081987
156 *
!!! 注意：数据及存储类型设置不当，如类型设置过小就会使得一些数据无法正常输入，甚至在两个变量或多个变量的观测值进行数学计算时结果出错。
157 * recast命令：更改数值型变量的存储类型
158 *操作实例->
159 sysuse auto, clear
160 list gear_ratio in 1/5
161 display gear_ratio
162 recast int gear_ratio, force //recast命令用于更改变量的存储类型
163 display gear_ratio
164 list gear_ratio in 1/5
165 * compress命令：精简资料的存储格式
166 *-2.3.5.定义变量的显示格式
167 * format命令：
168 * 语法1: format varlist %fmt / format %fmt varlist
169 *
*读懂%fmt的常见格式含义。如%-18s靠左列印于屏幕上；若%18s，则靠右列印；若%~18s,则居中列印。
170 * 语法2: format [varlist]
171 *操作实例->
172 help format
173 list price gear_ratio in 1/5
174 format price %6.1f
175 format gear_ratio %6.4f
176 list price gear_ratio in 1/5
177 *-2.3.6.数值型分类变量和字符变量的转换（分类变量要定义值标签）
178 * encode命令：将字符变量转换为分类数值变量。
179 * 语法: encode varname [if] [in], generate(newvar) [label (name) noextend]
180 *操作实例->
181 webuse hbp2, clear
182 describe sex
183 encode sex, generate(gender) label(sex1b1)
184 describe gender
185 * decode命令：将分类数值变量转换为字符变量。注意：无值标签的数值变量不适用。
186 * 语法: decode varname [if] [in], generate(newvar) [maxlength(#)]
187 *操作实例->
188 webuse hbp2, clear
189 describe sex
190 label define gender 1 "female" 2 "male"
191 replace sex = "other" in 2
192 encode sex, generate(gender)
193 label list gender
194 decode gender, generate(sex2)
195 tostring gender, generate(sex3)
196 *-2.3.7.包含数值数据的字符型变量与数值型变量转换
197 * destring命令：包含数值数据的字符型变量转换为数值型变量

```

```

198 *      语法: destring [varlist], [generate (newvarlist) | replace] [options]
199 *      常用参数: ignore ("chars") 删除字符变量中的非数值字符
200 *      force将非数值字符转换为缺失值
201 *      replace: 转换后的变量值替代原变量值
202 *操作实例->
203 webuse destring2, clear
204 describe date
205 list date
206 destring date, ignore(" ") replace
207 destring price, generate(p1) ignore("$ ,")
208 encode price, generate(p2)
209 *      tostring命令: 将数值变量转换为字符变量
210 *      语法: tostring varlist, [generate (newvarlist) | replace]
211 *操作实例->
212 webuse tostring, clear
213 describe
214 list
215 tostring year day, generate(y1 d1)
216 decode year, generate(y2)
217 *-2.3.8.新变量生成
218 *      generate/egen命令:
219 *      语法1: generate [type] newvar=exp [if] [in]
220 *      语法2: egen [type] newvar=fcn(arguments) [if] [in] [, options]
221 *操作实例->
222 sysuse auto, clear
223 gen id=.
224 gen lowprice=1 if price<4500
225 replace lowprice=0 if lowprice==.
226 gen lowprice2=(price<4500)
//注意: 如果price中有缺失值时, 该命令会将price为缺失值的样本的lowprice赋值为0, 因为stata会将缺失值
视为最大正值。
227 webuse egenxmpl2, clear
228 by dcode, sort: egen medstay = median(los)
229 gen slos1=sum(los)
230 egen slos2=sum(los)
231 webuse egenxmpl3, clear
232 egen byte differ = diff(inc1 inc2 inc3)
233 webuse egenxmpl4, clear
234 egen hsum = rowtotal(a b c)
235 egen hnonmiss = rownonmiss(a b c)
236 egen hsd = rowsd(a b c)
237 *      语法3: generate newvar=recode(varname, num1,num2, num3, ..., numk)
238 *操作实例->
239 gen priceg=recode(price, 2000, 4000, 6000)
//在新变量priceg中, price<=2000的样本赋值为2000, >2000&lt;=4000 赋值为4000, >4000赋值为6000。
240 *      recode命令: 分类变量再编码
241 *      语法4: recode varlist (rule) [(rule) ...] [, generate(newvar)]
242 *操作实例->
243 webuse fullauto, clear
244 recode rep77 rep78 (1 2 = 1 "Below average") (3 = 2 Average) (4 5 = 3 "Above average"), pre(new)
label(newrep)
245 label list repair newrep
246 *      replace命令: 变量值的修改
247 *      语法5: replace oldvar=exp [if] [in] [, nopromote]
248 *操作实例->
249 webuse genxmpl2, clear
250 generate lastname=word(name,2)
251 list
252 replace lastname=word(name,1)
253 list
254 *-2.3.9.缺失值的处理
255 *      在调查中, 经常用88,
99,888,999,... 等来表示不知道或不清楚。而这些信息在数据处理时, 均属于非有效信息, 是无法进入统计分析的。
因此, 这些都可视为观测缺失值。Stata中通常是用"."来表示这类观测缺失值。
256 *      mvencode命令: 用特定含义的数值来表示观测缺失值。
257 *      语法: mvencode varlist [if] [in], mv(#|mvc=# [\ mvc=#...] [\ else=#]) [override]
258 *      mvdecode命令: 将特定含义的数值处理为缺失值
259 *      语法: mvdecode varlist [if] [in], mv(numlist | numlist=mvc [\ numlist=mvc...])

```

```

260 *操作实例->
261 sysuse auto,clear
262 list rep78 foreign if rep78 == .
263 mvencode _all, mv(999) override
264 mvencode rep78 if foreign == 0, mv(998)
265 mvdecode rep78, mv(998=. \ 999=.a)
266 *-2.3.10.变量（观测值）的剔除、保留和显示
267 * drop/keep命令:
268 * 语法1: drop/keep varlist //变量
269 * 语法2: drop/keep if exp //观测值
270 *操作实例->
271 sysuse census
272 drop pop*
273 drop if medage > 32
274 * list命令:
275 * 语法: list varlist [if] [in] [,options]
276 *-2.3.11.数据的排序
277 * sort命令/gsort命令
278 * 语法1: sort varlist [in]
279 * 语法2: gsort [-] varname [[-] varname ...] [, generate (newvar)
280 * 语法3: gsort [-] varname [[-] varname ...] [, generate (newvar) mfirst], 可同时升、降序排序。参数mfirst表示将缺失值放在前面。
281 *操作实例->
282 sysuse auto
283 keep make mpg weight
284 sort make mpg weight, stable
285 list in 1/10
286 gsort + make -mpg - weight
287 list in 1/10
288 *-----
289 *->2.4数据集的合并和附加
290 *-----
291 *-2.4.1.数据集的合并
292 * merge命令: 1对1匹配合并, 1对多（多对1）匹配合并, 多对多匹配合并, 按观测值1对1匹配合并
293 * 语法1: merge 1:1 varlist using filename [, options] //
294 * 语法2: merge m:1 varlist using filename [, options]
295 * 语法3: merge 1:m varlist using filename [, options]
296 *操作实例->
297 webuse overlap2, clear
298 merge 1:m id using https://www.stata-press.com/data/r16/overlap1, update replace
299 list
300 *-2.4.2.数据集的附加
301 * append命令: append using filename [filename ...] [, options]
302 *操作实例->
303 sysuse auto, clear
304 keep if foreign == 0
305 save domestic
306 sysuse auto, clear
307 keep if foreign == 1
308 keep make price mpg rep78 foreign
309 append using domestic, keep(make price mpg rep78 foreign) generate(app1)
310 *=====over=====
311

```