

# 统计学 (Stata实现)

## 第一次上机课

赵震宇 (2120253538)

南开大学 国际经济研究所

zzynankai@outlook.com







xishanyu2.github.io

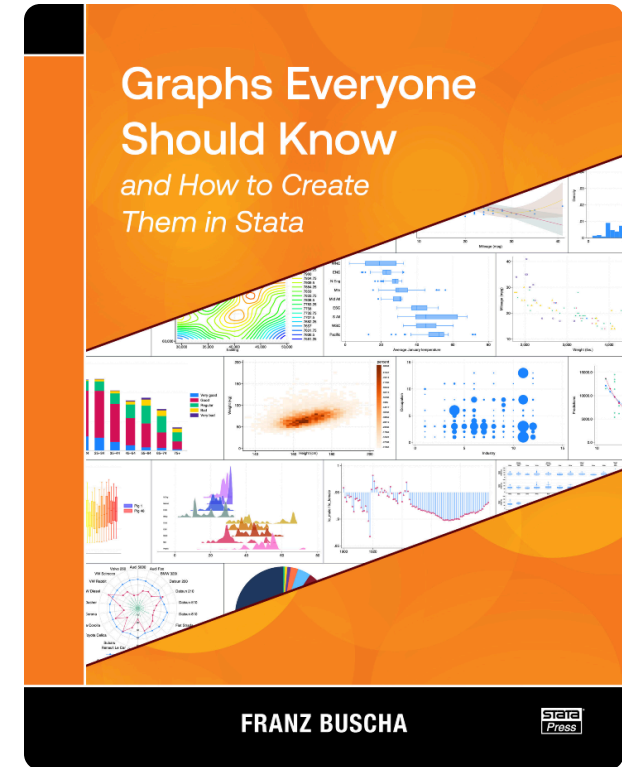
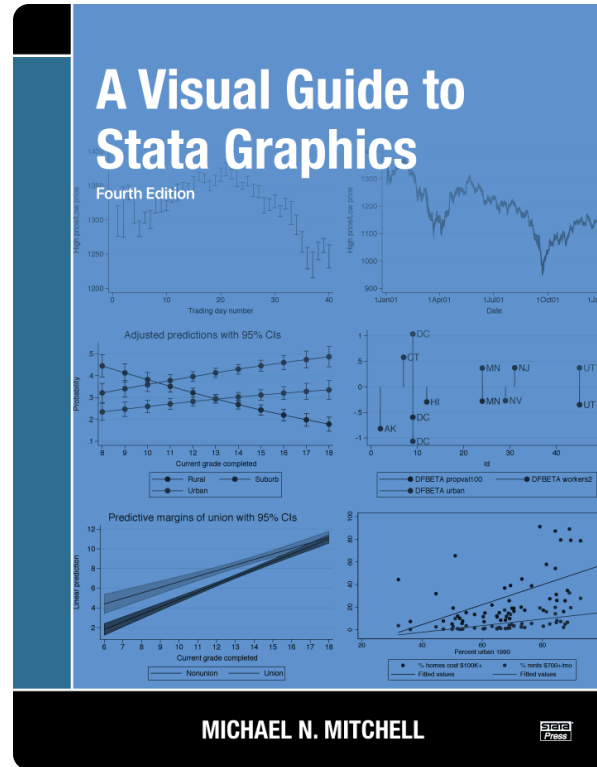
2025 年 10 月 16 日

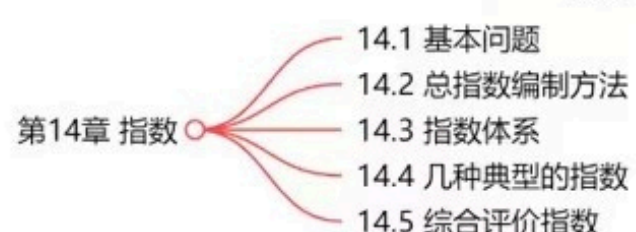
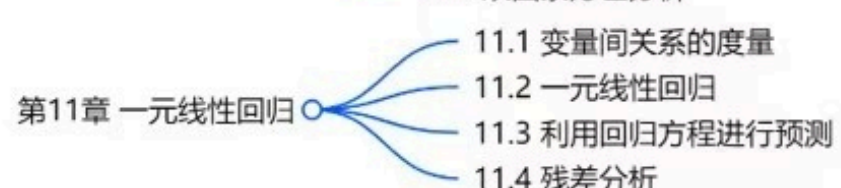
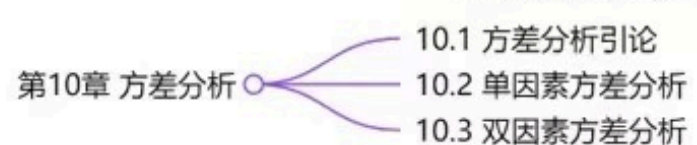
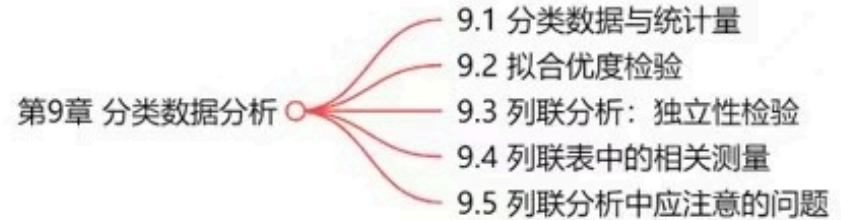
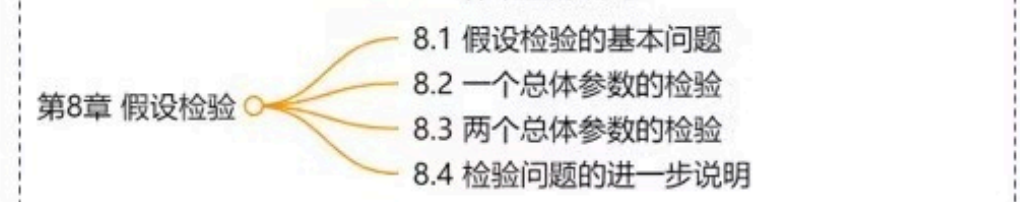
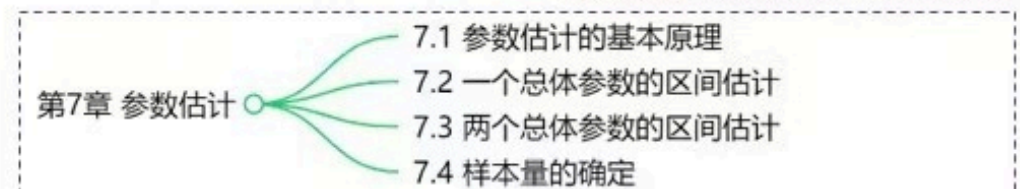
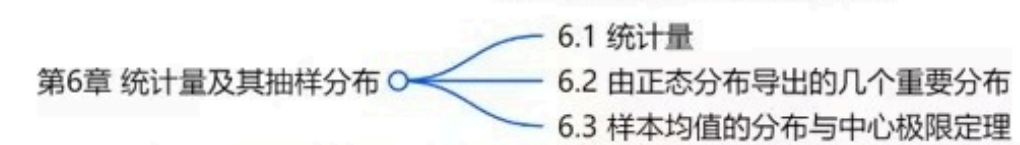
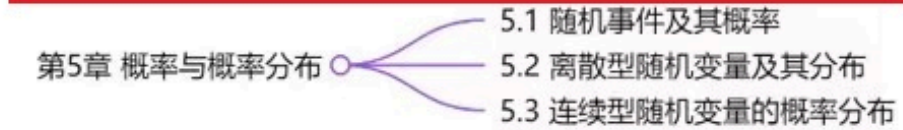
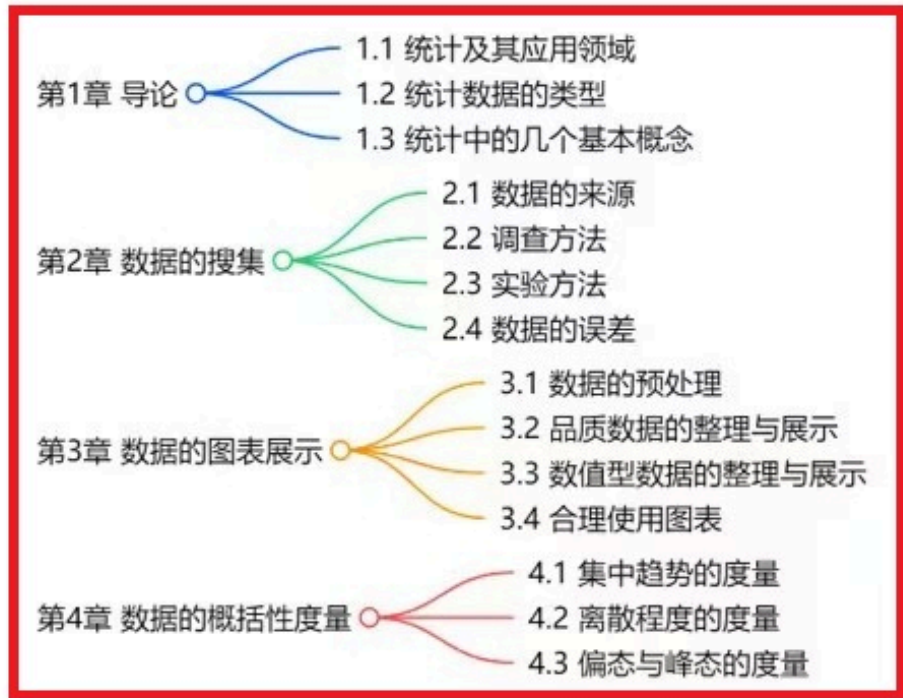
## 参考资料

- 张甜, 杨维忠编著. *Stata统计学与案例应用精解*. 清华大学出版社. 2025
- 谢慧明. *经济学实证论文写作讲义: 方法与应用*. 经济科学出版社. 2024
- 陈强. *计量经济学及Stata应用*. 高等教育出版社. 2023
- 陈舒艳编. *统计学: Stata应用与分析*. 机械工业出版社. 2020
- 连玉君老师<sup>↑</sup>、司继春老师<sup>↑</sup> 课件、资料等



- Stata Conferences  :
  - Data visualization with Stata 
- Stata Bookstore  :
  - A Visual Guide to Stata Graphics, Fourth Edition 
  - Graphs Everyone Should Know and How to Create Them in Stata 
- Documentation | Stata 





# CONTENTS

## 目录

1 Stata基础 (略) ↑

2 数据可视化 ↑

3 描述性统计 ↑

4 课后习题 ↑

# 1.引言

## 1.1 Why use Stata?

1. Excel, SPSS, EViews: 点击式 v. s Stata, R, Python: 命令式——易于复现!
  2. 封装好的命令（简洁!）、丰富的内置数据集与外部命令
  3. 全新的画图风格（从Stata18开始）
  4. 后续学习：计量经济学、因果推断、机器学习（从Stata19开始）、经济学实证论文写作
- .....

## 1.2 Why visualize data?

“ A picture is worth a thousand words. — Barnard (1927)

## 1.3 Why use Stata for data visualization?

- Creating graphs in Stata is easy.
- Stata supports a wide variety of plots.
- Stata graphic commands are highly customizable and extensible.

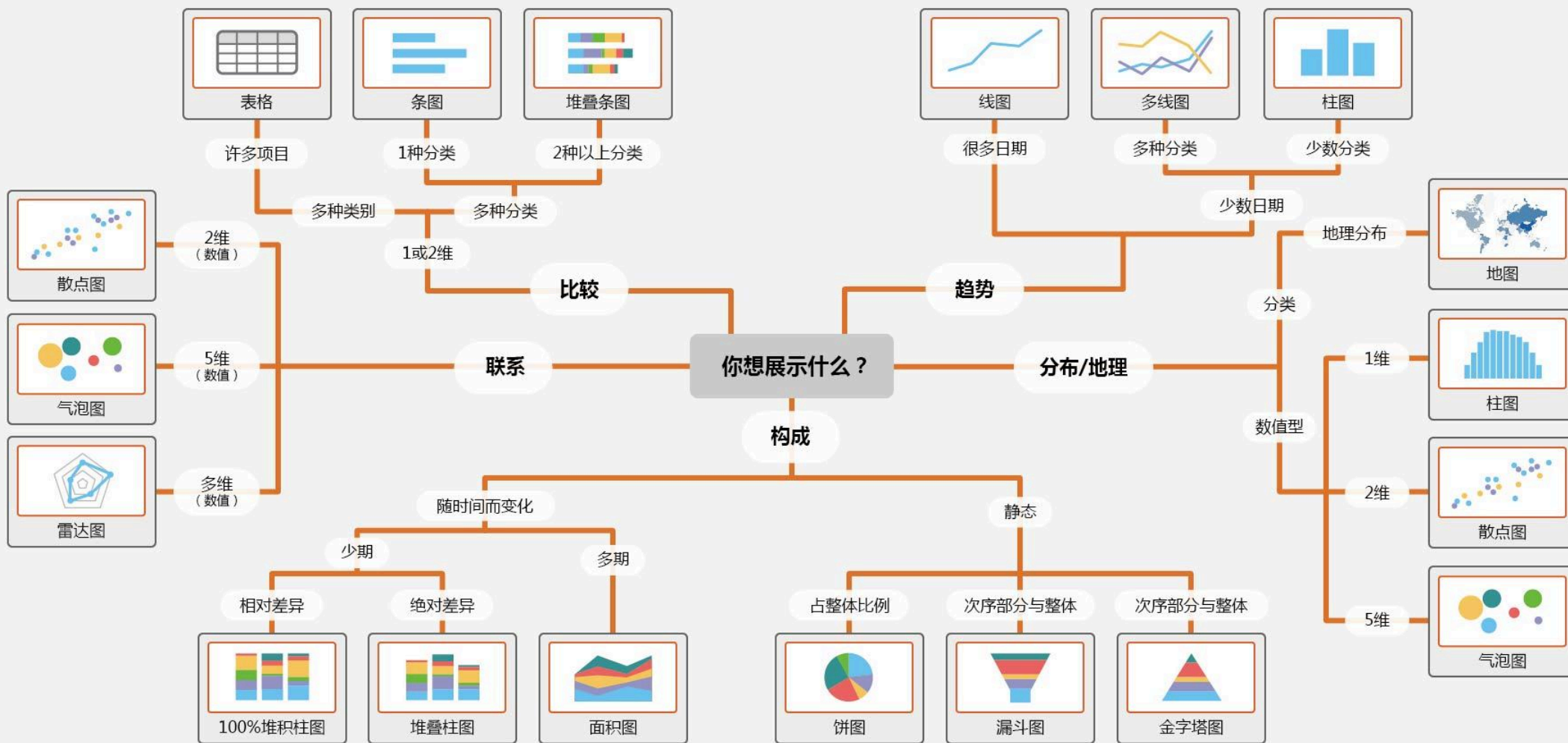
## 2. 数据可视化

## 2.0 数据图表简史 [video1 ↗](#) [video2 ↗](#)

“ Graphics that help us think faster or see a book's worth of information on a single page are the key to unlocking new discoveries.

“ Graphicacy enables us to harness our built-in GPU to process mountains of data and find the veins of gold hiding within.

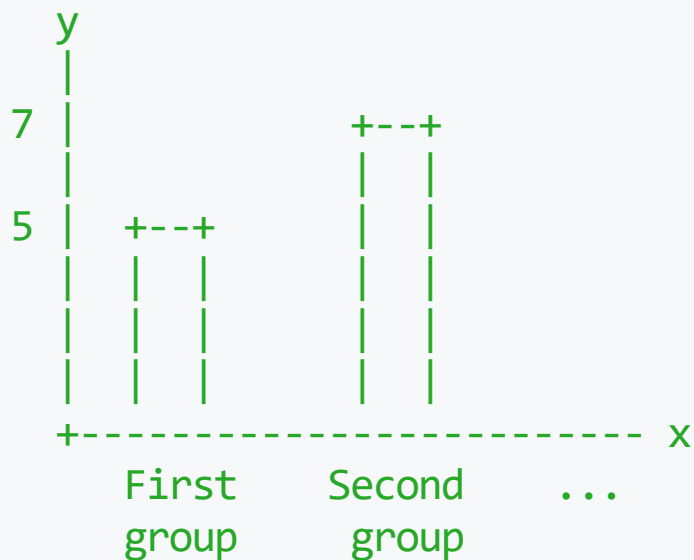
### 如何选择图表的类型？



## 2.1 类别数据可视化

### 条形图

```
graph bar (mean) numeric_var, over(cat_var)
/*
```



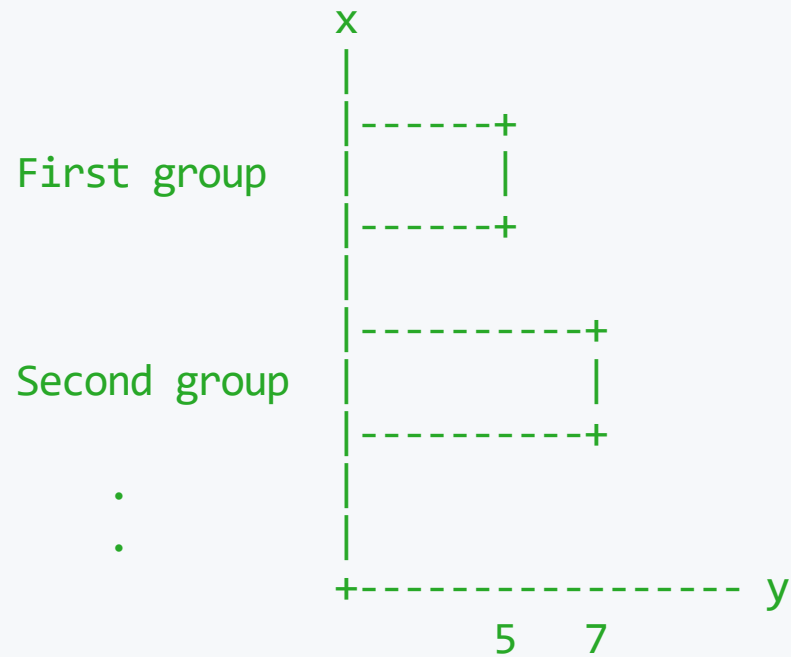
numeric\_var must be numeric;  
statistics of it are shown on  
the y axis.

cat\_var may be numeric or string;  
it is shown on the categorical  
x axis.

```
*/
```

```
graph hbar (mean) numeric_var, over(cat_var)
```

```
/*
```



same conceptual layout:  
numeric\_var still appears  
on y, cat\_var on x

```
*/
```

```
graph bar (mean) y, over(division) over(region)
/*
      +--+  +--+  +--+          +--+          +--+
      |||   |||   |||         |||           |||
      |||   |||   |||         |||           |||
-----
div_1  div_2  div_3      div_1  div_2  div_3
      region_1                region_2
*/
```

\* Simple bar chart

```
graph bar (count), over(x)
```

\* Horizontal percentage bar chart

```
graph hbar (percent), over(x)
```

\* A sorted bar chart

```
graph hbar (percent), over(x, sort(1) descending)
```

\* Multiple bar charts

```
graph bar (count), over(x) by(y)
```

\* Multiple bar charts with two or three over() options

```
graph bar (count), over(x) over(y) over(z)
```

## 饼图

\* Basic pie chart

```
graph pie, over(x)
```

\* Labeling pies

```
graph pie, over(x) plabel(_all percent)
```

\* Exploding pies

```
graph pie, over(x) plabel(_all name) pie(2, explode(10))
```

\* Multiple pie charts with by()

```
graph pie, over(x) by(y)
```

## 2.2 数据分布可视化

### 直方图

\* A basic histogram

```
histogram x
```

\* Options that change bin sizes

```
histogram x, bin(10) name(g1, replace)
```

```
histogram x, bin(50) name(g2, replace)
```

```
histogram x, discrete name(g3, replace)
```

```
histogram x, width(10) name(g4, replace)
```

```
graph combine g1 g3 g2 g4
```

\* Options that change the y axis

```
histogram x, density name(g1, replace)
histogram x, fraction name(g2, replace)
histogram x, frequency name(g3, replace)
histogram x, percent name(g4, replace)
graph combine g1 g3 g2 g4
```

\* Adding kernel density plot

```
histogram x, width(5) kdensity kdenopts(bwidth(5))
```

\* Multiple histograms with by()

```
histogram x, by(y)
```

\* Multiple two-way histograms

```
twoway (histogram x if y == 1, color(edkblue%30)) ///
(histogram x if y == 2, color(erose%30)) ///
(histogram x if y == 3, color(stone%30)) ///
(histogram x if y == 4, color(emerald%30))
```

## 核密度估计图

\* A basic kernel density plot

```
kdensity x
```

\* Basic visualization options

```
kdensity x, recast(area) fcolor(eltblue) lcolor(black)
```

\* Density plots with different bandwidths

```
kdensity x, bwidth(10)
```

\* Multiple kernel density plots

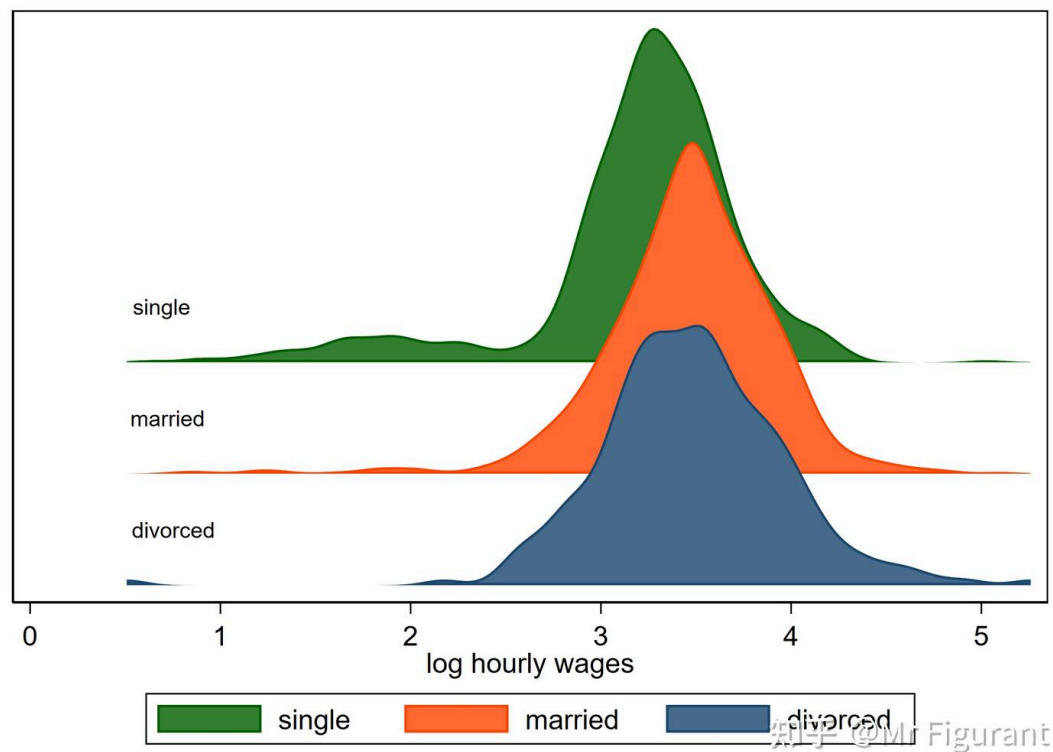
```
twoway (kdensity x, by(y))
```

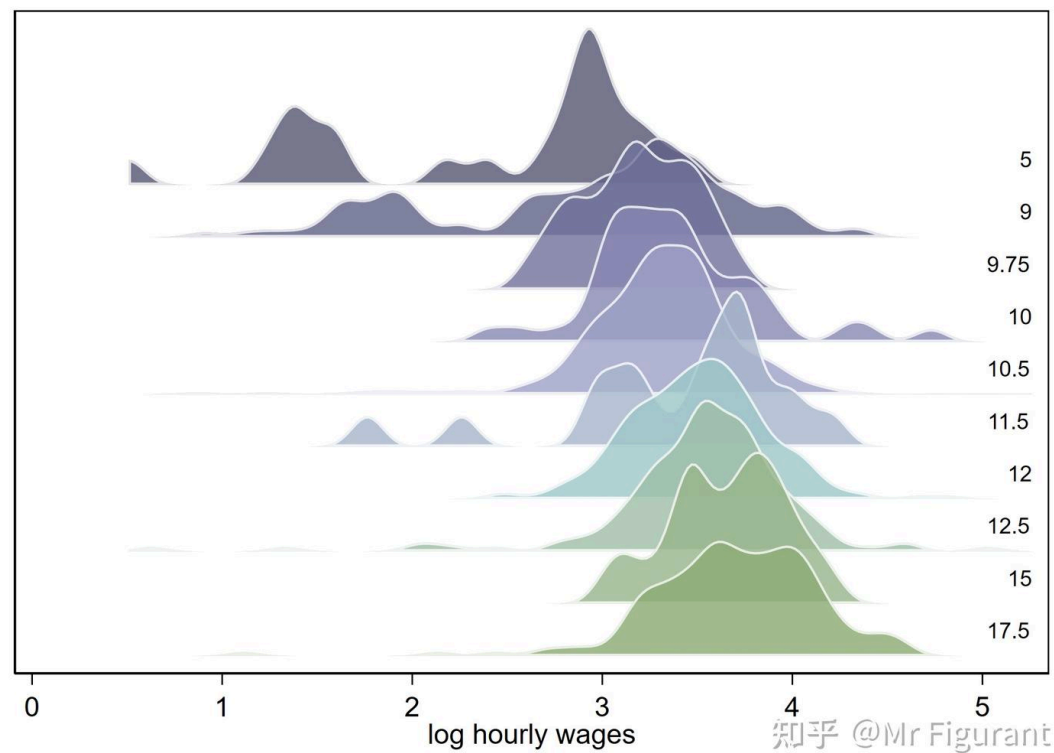
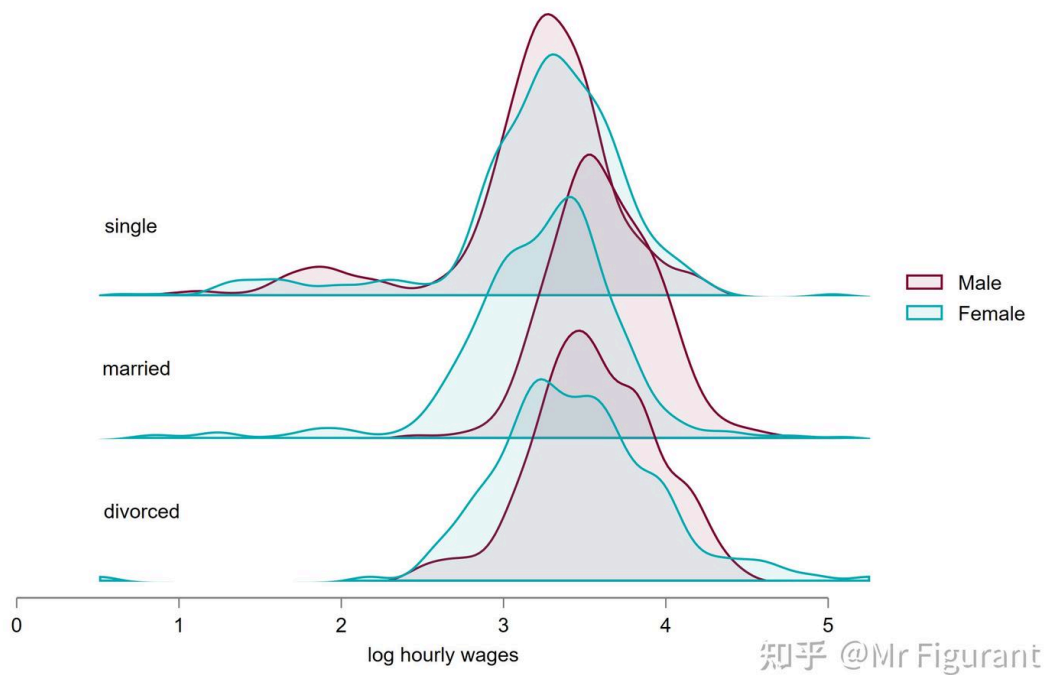
\* Multiple kdensity plots with if qualifiers

```
twoway (kdensity x if y == 1) (kdensity x if y == 2)
```

## 拓展阅读：joy\_plot

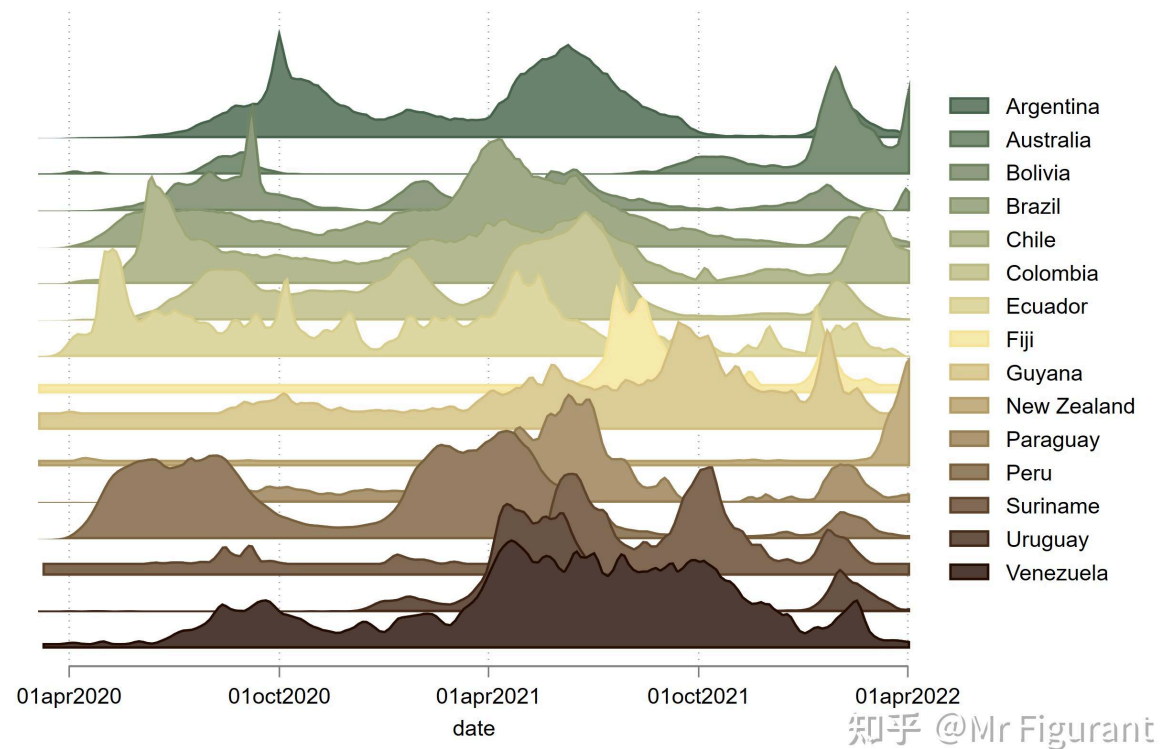
- Stata学习：如何绘制Fancy的核密度图？ - 知乎 [↗](#)

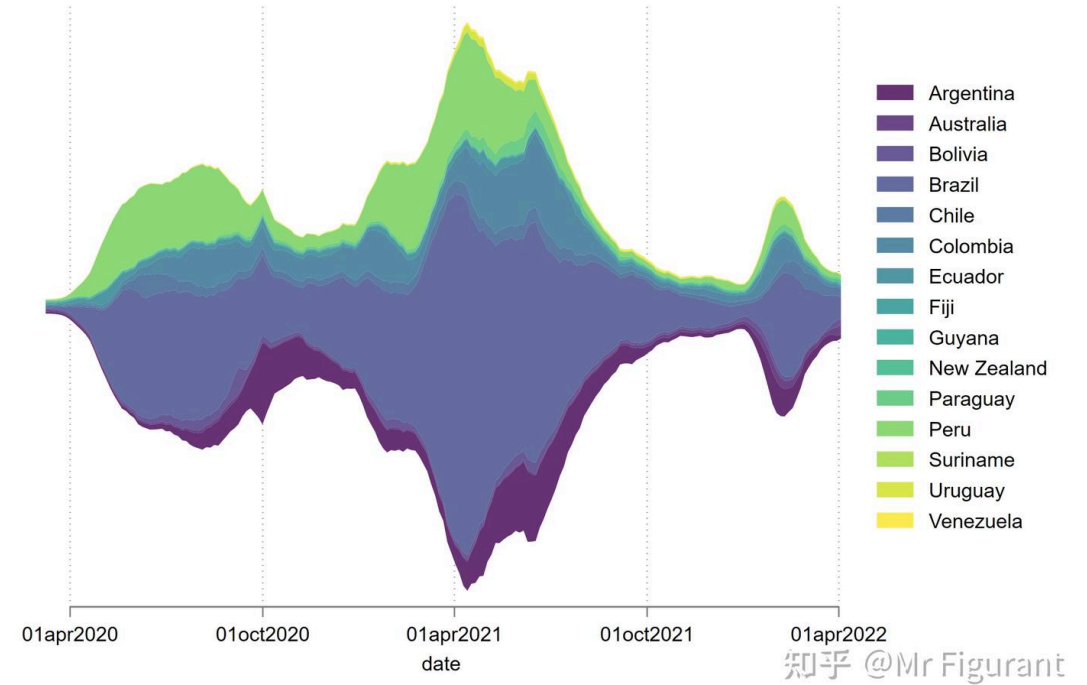
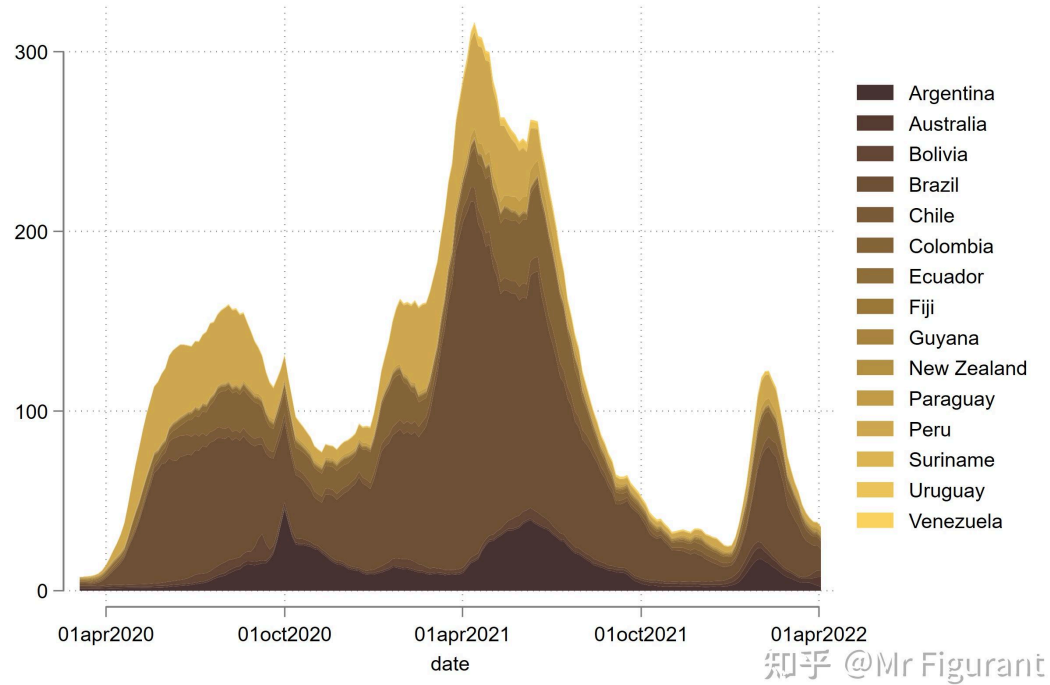




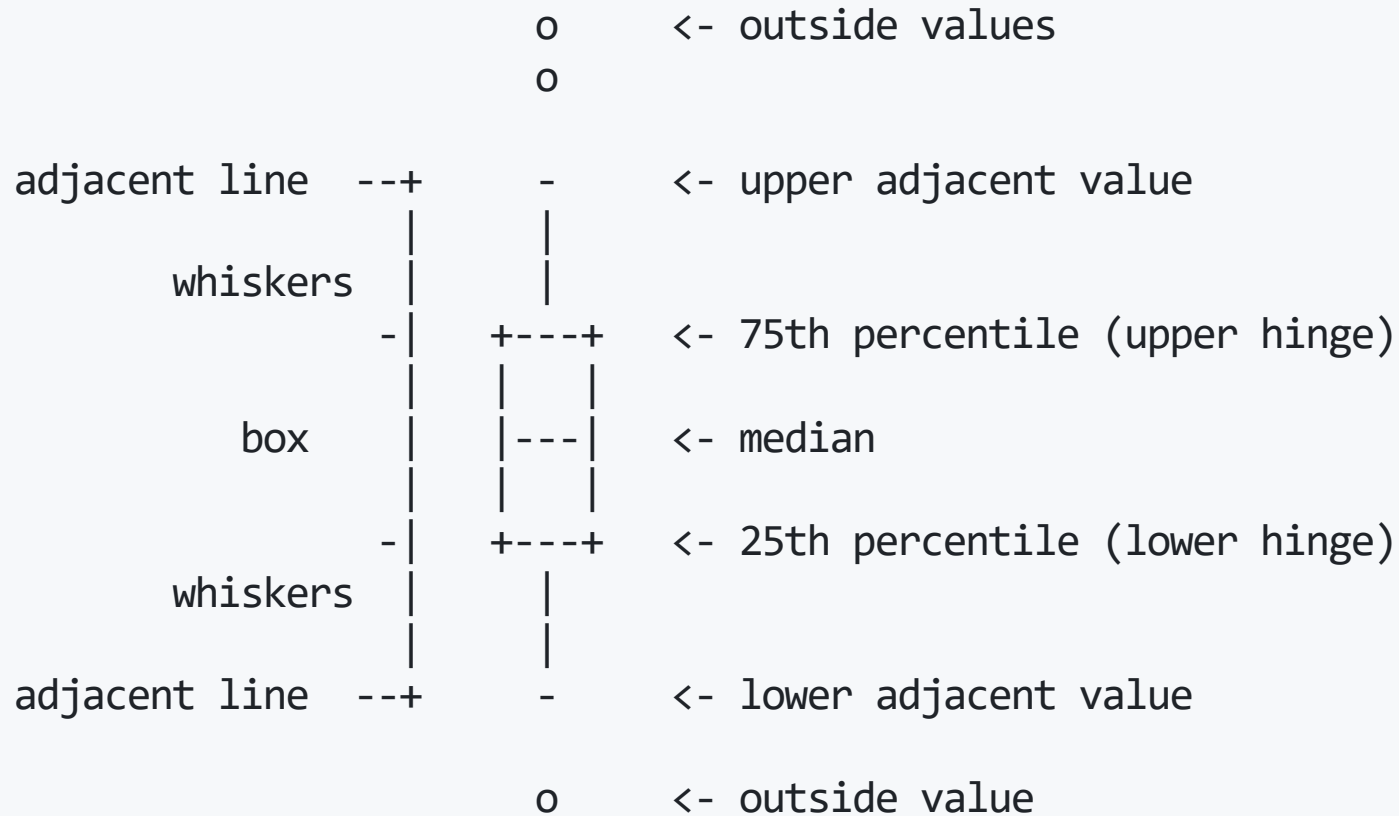
## 拓展阅读：脊线图

- Stata学习：如何绘制脊线图 ridgeline? - 知乎 [↑](#)



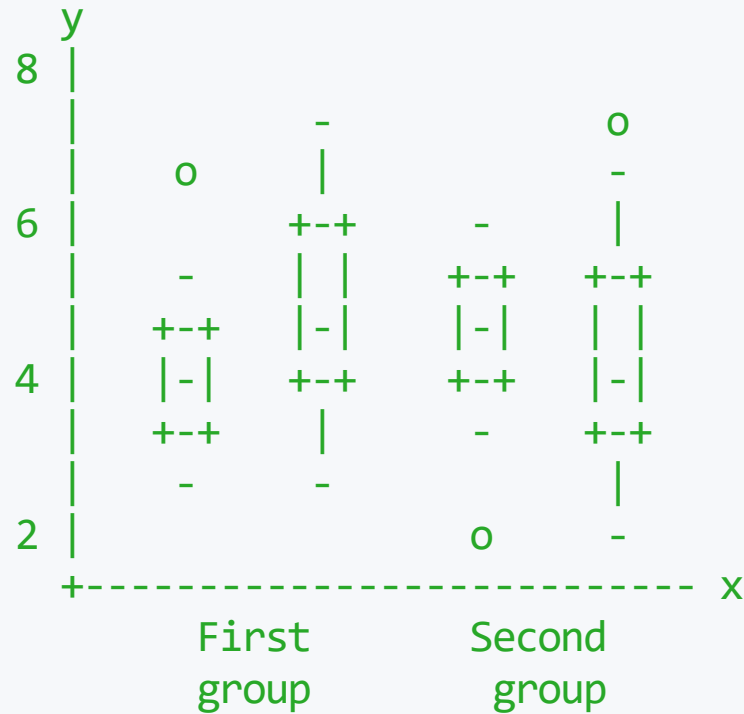


## 箱线图



```
graph box y1 y2, over(cat_var)
```

```
/*
```

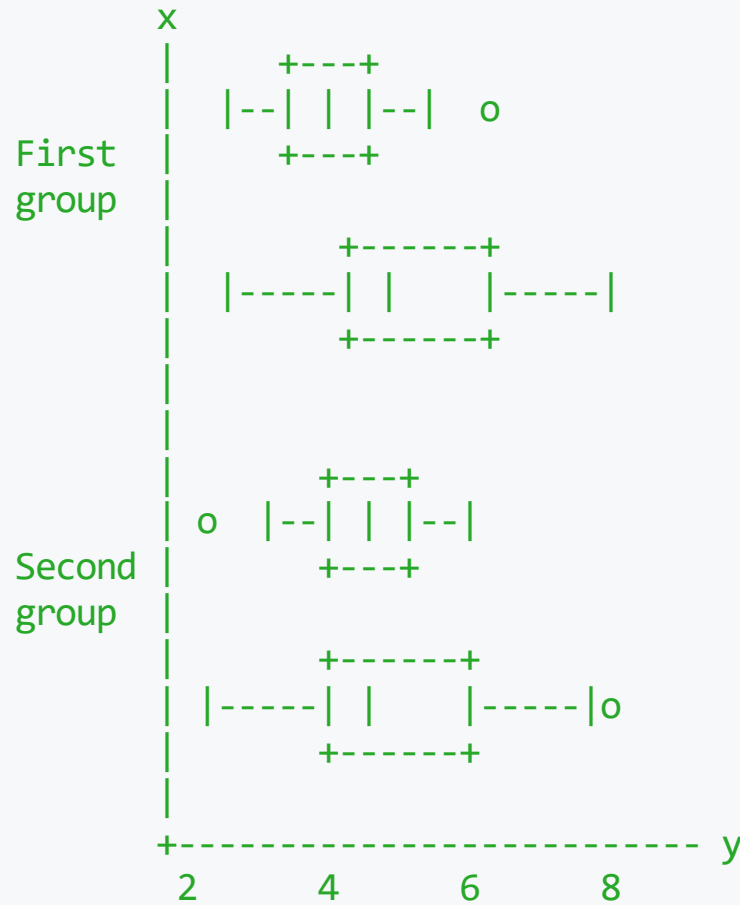


y1, y2 must be numeric;  
statistics are shown on  
the y axis

cat\_var may be numeric  
or string; it is shown  
on categorical x axis

```
*/
```

```
graph hbox y1 y2, over(cat_var)
/*
```



same conceptual layout  
as for graph box:

y1, y2 appear on y axis

cat\_var appears on x axis

```
* Basic box plot
graph box x

* Basic box plot with outliers
graph box x if y == 1

* Horizontal box plot
graph hbox x

* Histogram with box plot
graph hbox x, name(g1, replace)
histogram y, name(g2, replace) ytitle("")
graph combine g1 g2, cols(1) xcommon
```

\* Multiple box plots

```
graph box x y
```

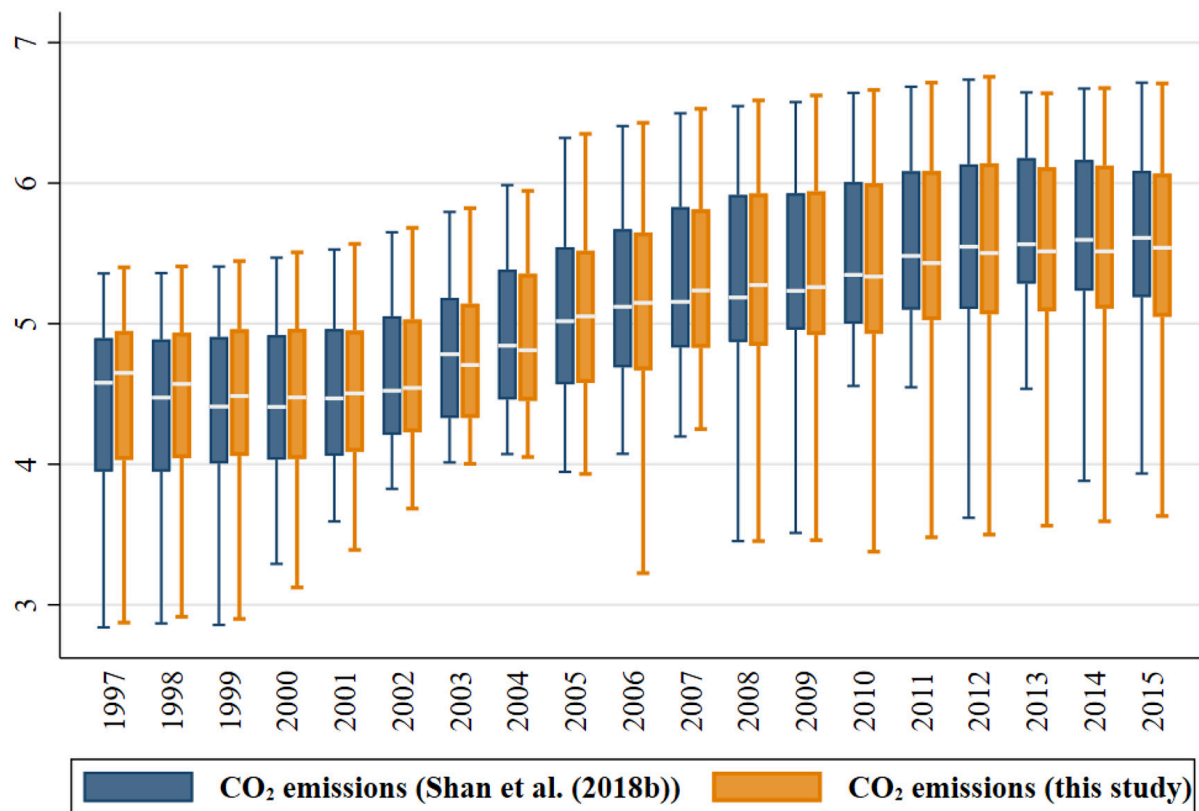
\* Multiple box plots over a categorical variable

```
graph box x y, over(z)
```

\* Ordering boxes

```
graph hbox x, over(z, sort(1))
```

## 论文阅读



Chen Y, Shao S, Fan M, et al. One man's loss is another's gain: does clean energy development reduce CO2 emissions in China? Evidence based on the spatial Durbin model[J]. Energy Economics, 2022, 107: 105852.

## 小提琴图

```
ssc install vioplot
```

```
* A violin plot
```

```
vioplot x
```

```
* Violin plot with kernel() option
```

```
vioplot x, kernel(gaussian) bwidth(10)
```

```
* A horizontal violin plot
```

```
vioplot x, horizontal
```

```
* Multiple violin plots
```

```
vioplot x y, over(z)
```

```
* Multiple violin plots with custom widths
```

```
vioplot x y, over(z) barwidth(500) dscale(200)
```

# 自学

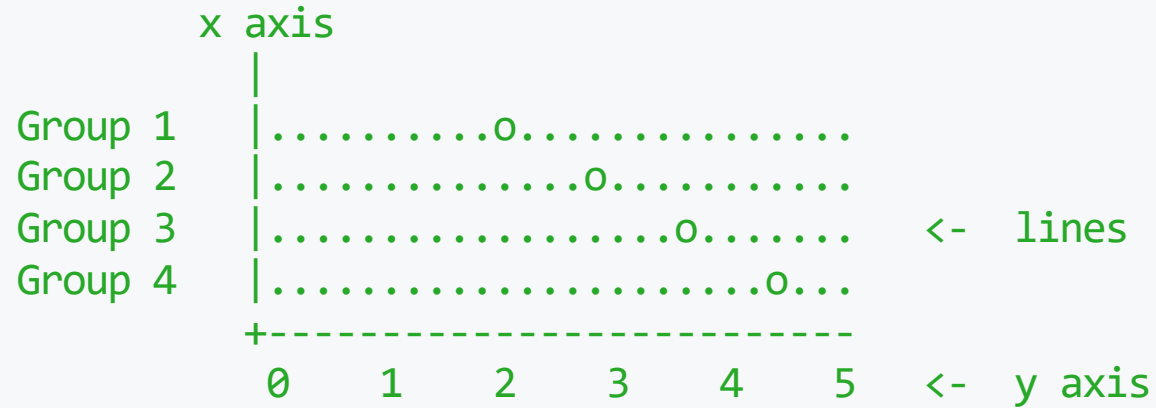
点阵图 `dotplot`

茎叶图 `stem`

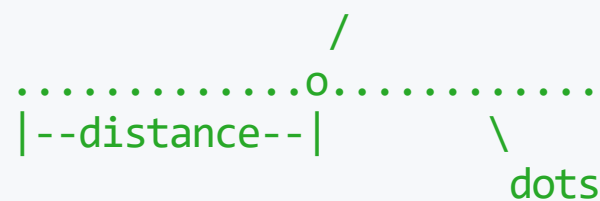
# 区分 graph dot 和 dotplot

```
graph dot (mean) numeric_var, over(cat_var)
```

```
/*
```



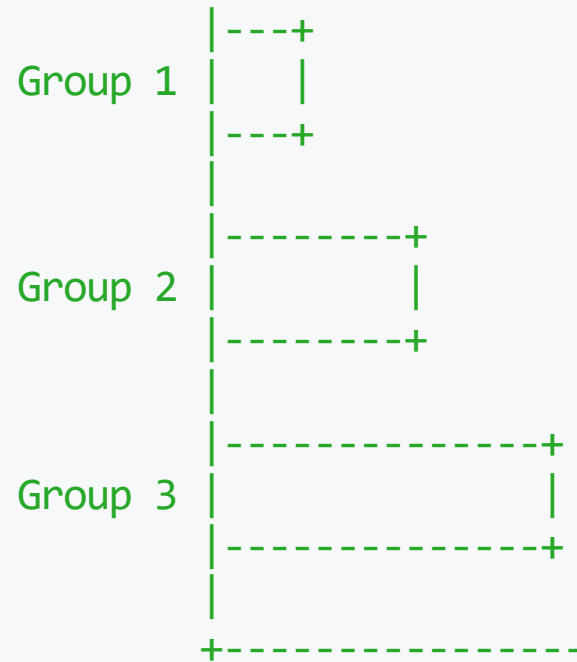
marker



```
*/
```

```
graph hbar y, over(group)
```

```
/*
```



```
*/
```

```
graph dot y, over(group)
```

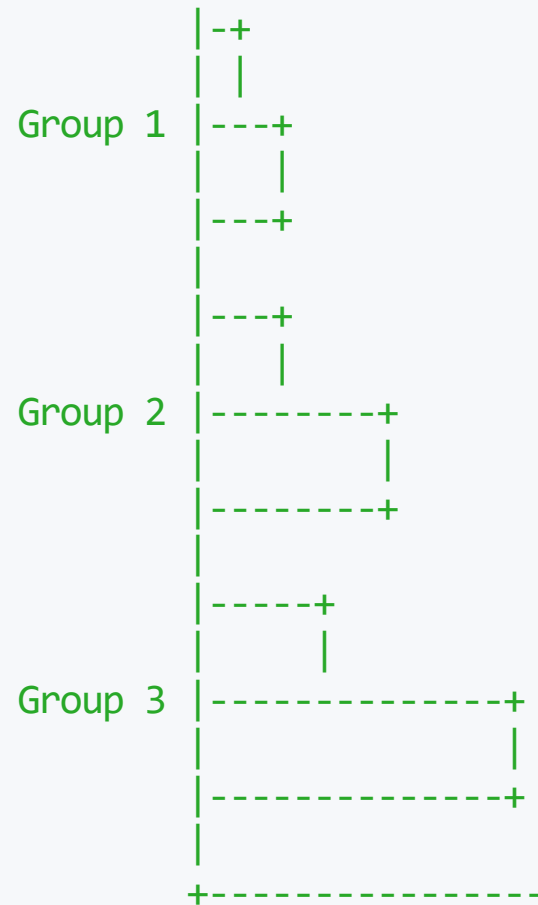
```
/*
```



```
*/
```

```
graph hbar y1 y2, over(group)
```

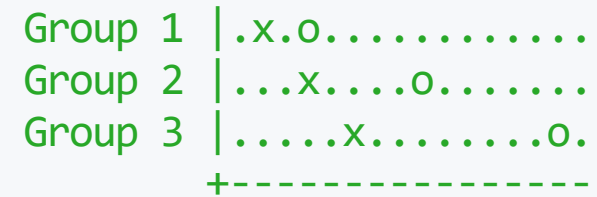
```
/*
```



```
*/
```

```
graph dot y1 y2, over(group)
```

```
/*
```



```
*/
```

## 2.3 变量间关系可视化

### 散点图

\* Basic scatterplot

```
twoway (scatter y x)
```

\* Square, colored, and sized markers

```
twoway (scatter y x in 1/15, msymbol(S) mcolor(edkblue%70) msize(*3) mlabel(make))
```

\* Multiple y variables against one x variable

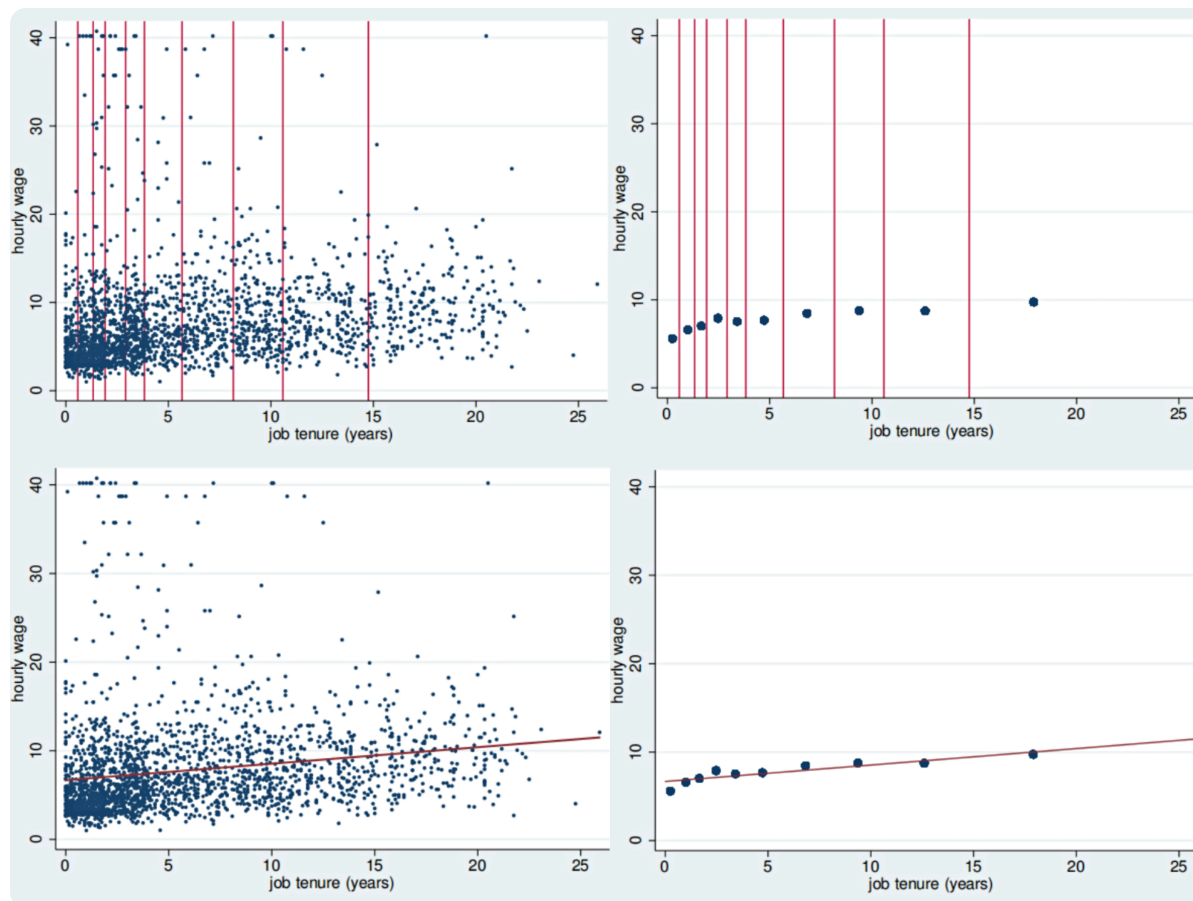
```
twoway (scatter x1 x2 x3 y, msymbol(S T D) mcolor(gs1 gs6 gs12))
```

\* Advanced syntax, two axes

```
twoway (scatter y x, xaxis(1) yaxis(1 2)) (scatter y z, xaxis(2))
```

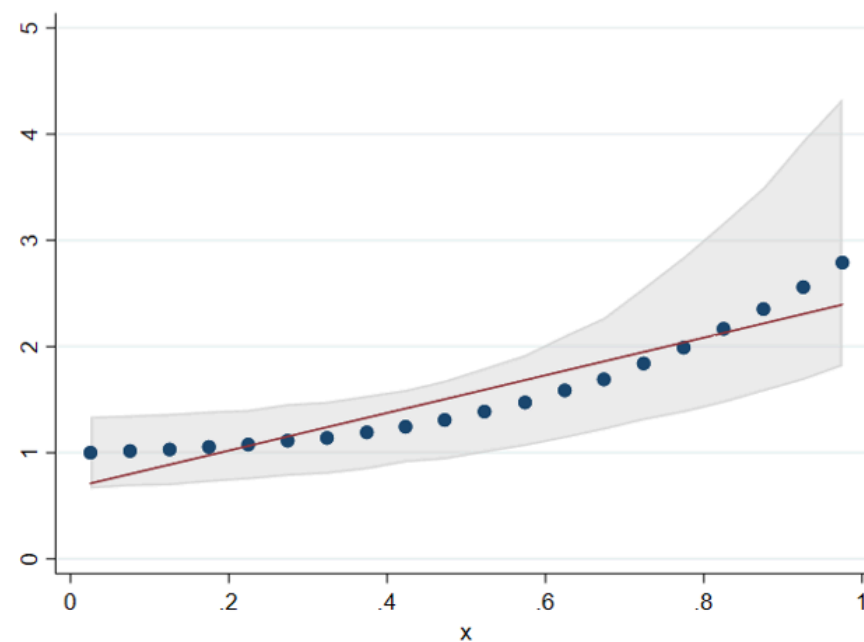
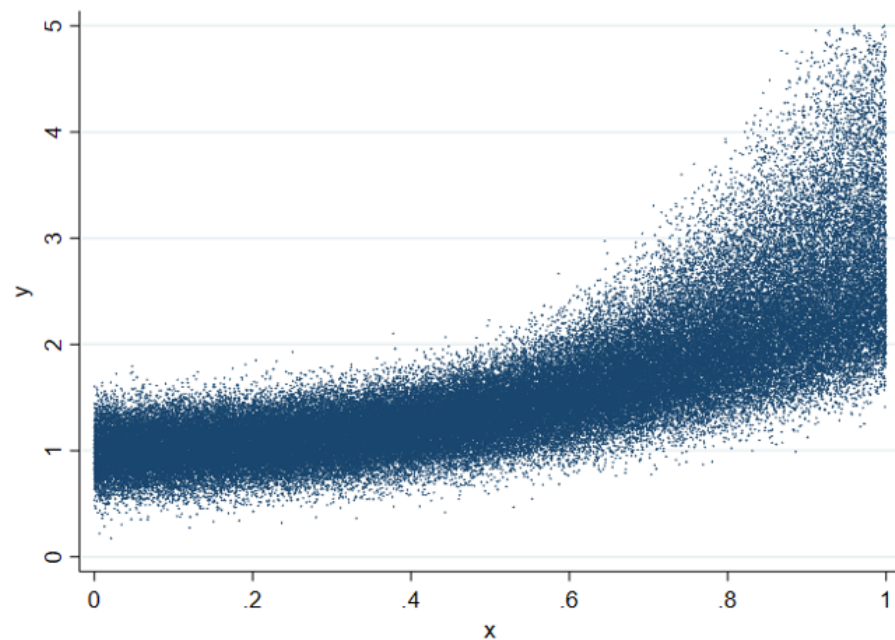
## 拓展阅读：binscatter

- Stata: 分仓散点图绘制-binscatter-binscatter2 ↑



## 拓展阅读：binscatter2

- Stata: 分仓散点图绘制-binscatter-binscatter2 ↗



## 相关系数矩阵与散点图矩阵

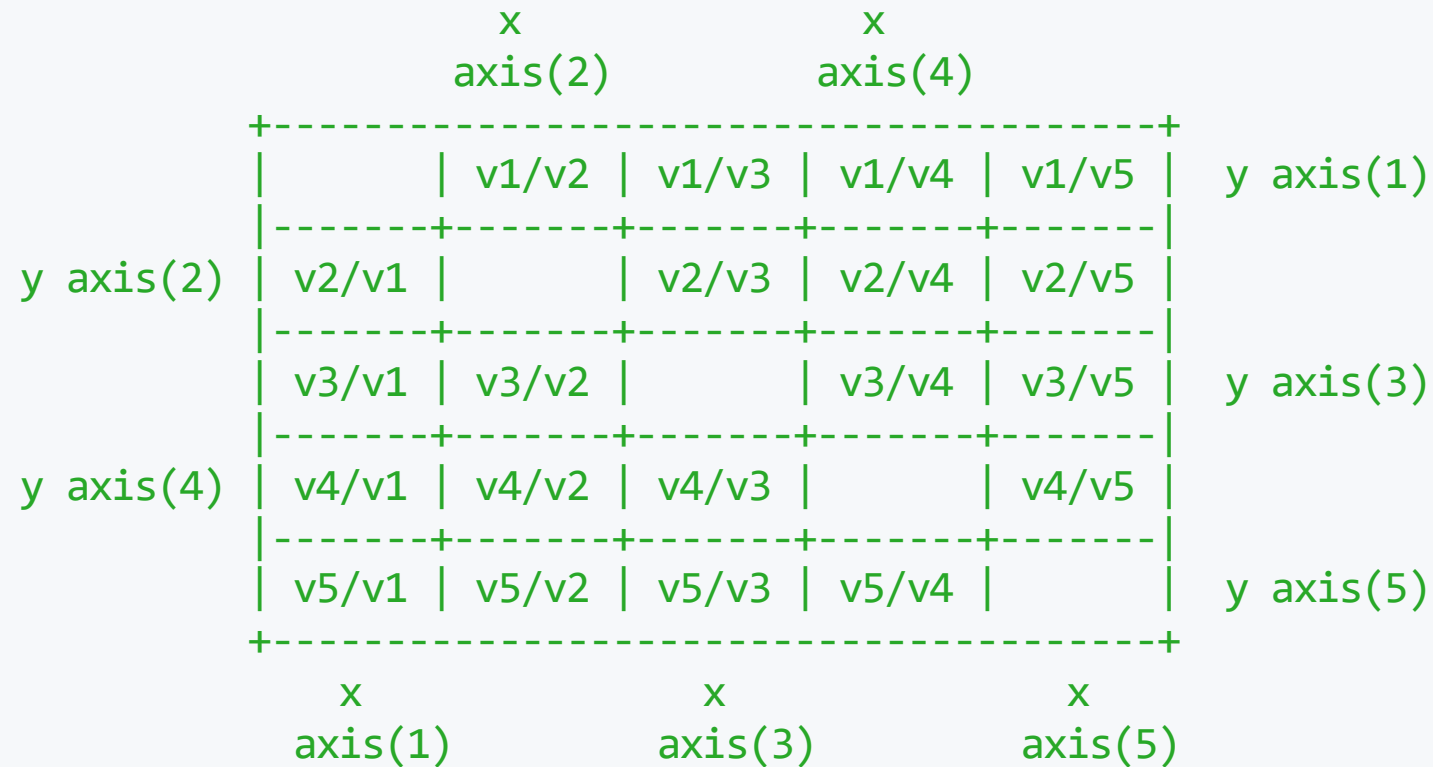
```
*Correlation matrix  
correlate v1 v2 v3 v4 v5
```

### 拓展：Pearson相关系数

如果样本中所有变量都没有缺失值，则corr命令和pwcorr命令输出的结果相同；否则，corr命令只针对所有变量都没有缺失值的子样本进行计算。而pwcorr则针对每一对变量进行计算，因此，每组变量在计算相关系数时使用的样本数可能是不同的。

```
graph matrix v1 v2 v3 v4 v5
```

```
/*
```



```
*/
```

```
graph matrix v1 v2 v3 v4 v5, half
```

```
/*
```

```

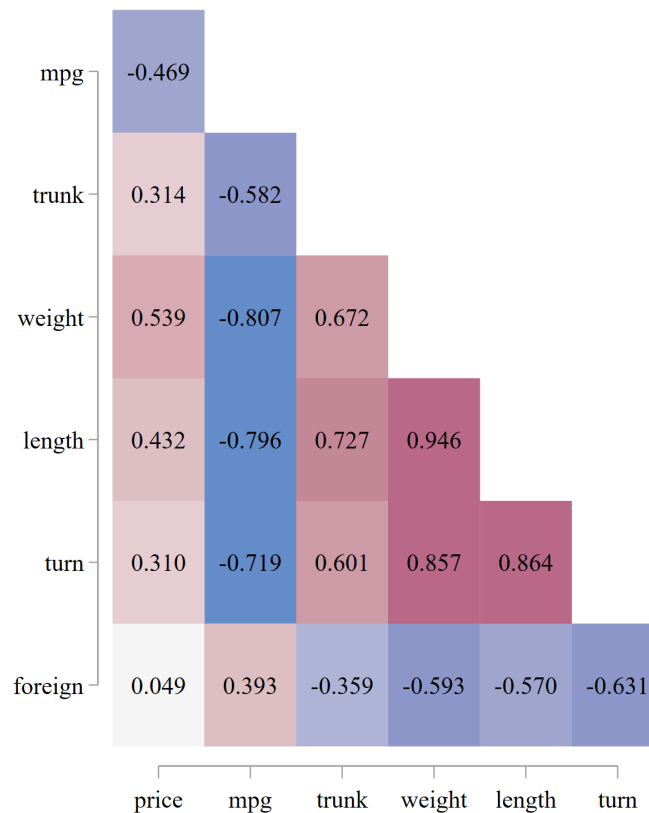
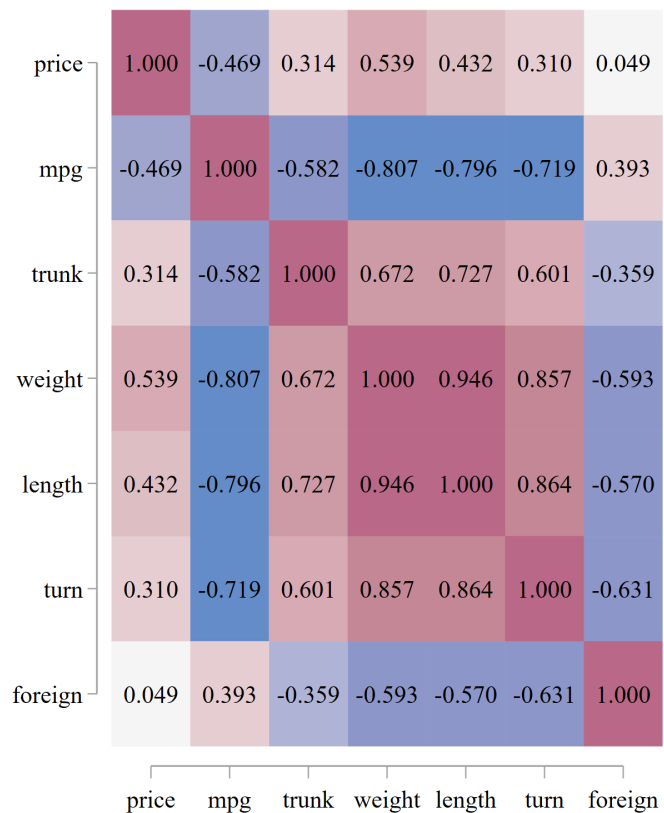
y axis(1) |          |
          |          |
          |-----+-----+
y axis(2) | v2/v1 |          |
          |-----+-----+
y axis(3) | v3/v1 | v3/v2 |          |
          |-----+-----+-----+
y axis(4) | v4/v1 | v4/v2 | v4/v3 |          |
          |-----+-----+-----+-----+
y axis(5) | v5/v1 | v5/v2 | v5/v3 | v5/v4 |          |
          |-----+-----+-----+-----+
          x         x         x         x         x
          axis(1) axis(2) axis(3) axis(4) axis(5)

```

```
*/
```

# 拓展：用热力图可视化相关系数矩阵

- Stata 可视化: heatplot-热力图 ↕



## 气泡图

```
* A basic bubble plot  
tway (scatter y x [fweight = z])
```

## 雷达图

```
ssc install radar  
  
* Radar chart with aspect ratio  
radar z x in 1/12, aspect(1) graphregion(color(white)) xlabel(, nogrid) //控制图形纵横比为1:1, 确保雷达图呈正圆形  
  
* Radar charts with multiple variables  
radar z x1 x2 x3, rlabel(0 10 20 30 40 50) aspect(1) legend(position(6) row(1)) xlabel(, nogrid) //legend图例  
  
* Radar charts with multiple variables and radial representation  
radar z x1 x2, rlabel(0 10 20 30 40 50) aspect(1) connected radial(x3) legend(position(6) row(1)) xlabel(, nogrid)
```

## 2.4 时间序列可视化

### 折线图

\* A basic line plot

```
line x year
```

\* Multiple line plots via line

```
line x1 x2 year, lcolor(gs12 gs4)
```

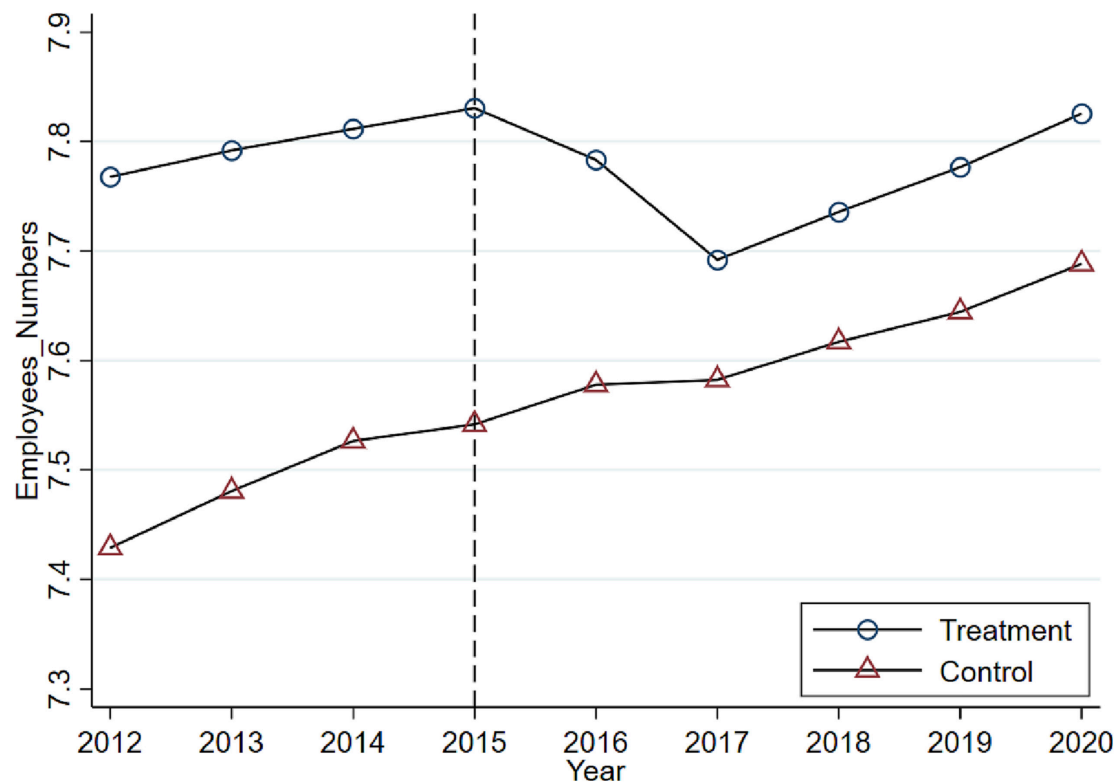
\* Multiple line plots via twoway

```
twoway (line x1 year, lcolor(gs12)) (line x2 year, lcolor(gs4))
```

\* Line plots with markers

```
twoway (line x year) (scatter x year, msymbol(D) msize(*0.5) mcolor(%50))
```

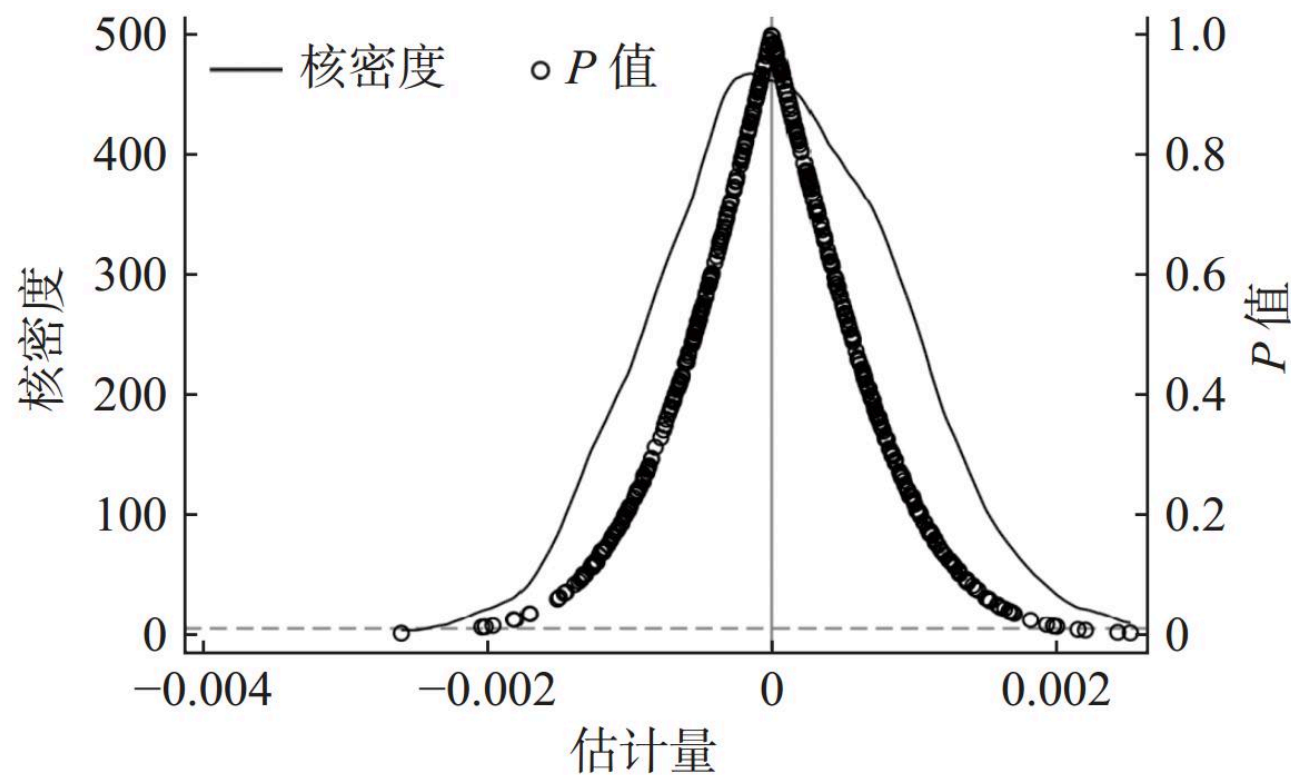
## 论文阅读



Liao T, Liu G, Liu Y, et al. Environmental regulation and corporate employment revisited: New quasi-natural experimental evidence from China's new environmental protection law[J]. Energy Economics, 2023, 124: 106802.

## 安慰剂检验——散点图+核密度图

- 平行趋势检验：DID的前提
- 安慰剂检验：DID是否稳健



## 2.5 图形叠加：twoway, addplot

- grunfeld.dta ↴

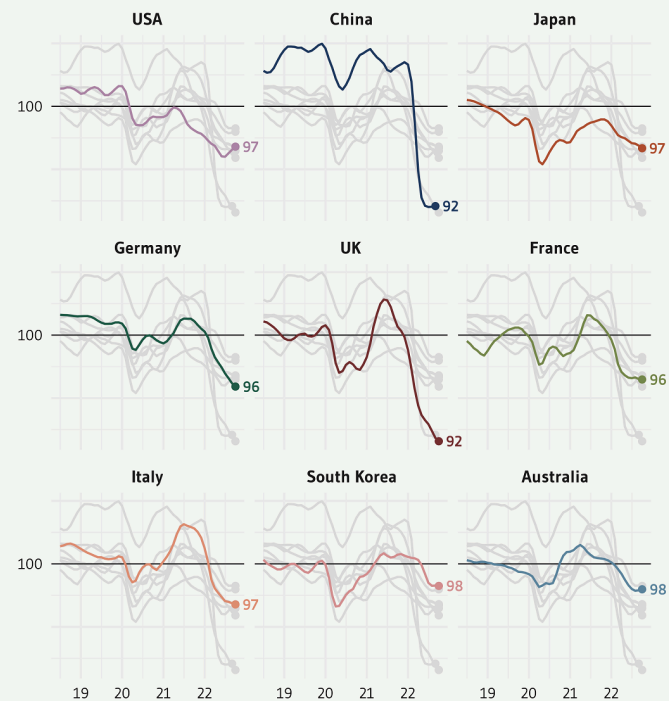
```
use "grunfeld.dta", clear

gen lnInv = ln(invest)
gen lnMV = ln(mvalue)

local y "lnInv"
local y "lnMV"
twoway line `y' time if company==3, ///
    lw(*2) lc(black) lp(solid)
foreach i of numlist 1 2 4/10{
    addplot: line `y' time if company==`i', ///
        lc(black*0.6%50) lp(solid) legend(off)
}
```

### Consumer Confidence Around the World

The consumer confidence indicator provides an indication of future developments of households' consumption and saving. An indicator above 100 signals a boost in the consumers' confidence towards the future economic situation. Values below 100 indicate a pessimistic attitude towards future developments in the economy, possibly resulting in a tendency to save more and consume less. During 2022, the consumer confidence indicators have declined in many major economies around the world.



Design: Gilbert Fontana  
Data: OECD, 2022

# 3. 描述性统计

```
sysuse auto
sum price
/*
  Variable |      Obs      Mean   Std. dev.   Min   Max
-----+-----
  price |      74  6165.257  2949.496   3291 15906
*/
```

```
sum price, detail
```

```
/*
```

```
Price
```

```
-----  
Percentiles      Smallest  
1%              3291          3291  
5%              3748          3299  
10%             3895          3667  
25%             4195          3748  
  
50%             5006.5  
  
75%             6342          13466  
90%             11385         13594  
95%             13466         14500  
99%             15906         15906  
*/  
  
Obs              74  
Sum of wgt.     74  
  
Mean             6165.257  
Std. dev.        2949.496  
  
Variance         8699526  
Skewness         1.653434  
Kurtosis         4.819188
```

## 3.1 描述水平的统计量

### 简单平均数

```
tabstat x, stats(mean) //tabstat x 默认是mean
```

### 加权平均数

```
gen xbar = m1 if x < x1 //m1为组中值  
replace xbar = m2 if x < x2 & x >= x1  
replace xbar = m3 if x < x3 & x >= x2  
replace xbar = m4 if x < x4 & x >= x3  
tabstat xbar, stats(mean)
```

## 中位数

```
tabstat x, stats(median)
```

## 四分位数

```
tabstat x, stats(q)
```

## 百分位数

```
tabstat x, stats(p10)
```

## 频数和众数

```
tabulate rep78
```

```
/*
```

Repair record 1978	Freq.	Percent	Cum.
1	2	2.90	2.90
2	8	11.59	14.49
3	30	43.48	57.97
4	18	26.09	84.06
5	11	15.94	100.00
Total	69	100.00	

```
*/
```

- N = 74?

## 3.2 描述差异的统计量

### 极差

```
tabstat x, stats(range)
```

### 四分位差

```
tabstat x, stats(iqr)
```

## 方差

```
tabstat x, stats(var)
```

## 标准差

```
tabstat x, stats(sd)
```

## 变异系数/离散系数

```
tabstat x, stats(cv)
```

## 异常值

```
sum x //返回r类命令  
gen z = (x - r(mean)) / r(sd)  
list x z if z > 3 //切比雪夫法则  
drop z if z > 3
```

### 拓展：r类命令、e类命令

e类命令：估计（estimation）命令，比如regress，通过e()调用

r类命令：其他命令，比如sum x，通过r()调用

## 3.3 描述分布形状的统计量

### 偏度系数

```
tabstat x, stats(skewness)
```

### 峰度系数

```
tabstat x, stats(kurtosis)
```

## 3.4 数据标准化

### 中心化

```
qui sum x  
gen c_x = x - r(mean)
```

### 标准化

```
qui sum x  
gen z_x = (x - r(mean)) / r(sd)
```

## 标准化2

```
egen z_x = std(x)
```

🔥 拓展：center（外部命令）

```
center x, prefix(c_) # 可加standardize, 简写至s
```

## 归一化

```
qui sum x  
gen n_x = (x - r(min)) / (r(max) - r(min))
```

## 3.5 正态性检验

### 直方图图示

```
hist x, norm
```

### 正态分位图

```
qnorm x, grid //add grid lines
```

## 正态性统计检验

```
sktest x
```

## 幂阶梯转换

```
ladder x  
qladder x //散点图矩阵  
gladder x //直方图矩阵
```

## 4. 课后习题

## 习题1 (数据可视化)

本题使用 `exercise3_3.csv` (请在 `pydata` 文件夹中寻找)

数据集来自R语言的 `mtcars`, 摘自1974年《美国汽车趋势》杂志, 包括32款汽车 (1973-1974年款) 的油耗、汽车设计和性能等共11个变量。根据该数据集绘制以下图形。

- (1) 绘制气缸数量 (`cyl`) 的直方图, 并为直方图叠加核密度曲线。
- (2) 绘制每加仑油行驶的英里数 (`mpg`) 和汽车自重 (`wt`) 两个变量的箱线图和小提琴图。
- (3) 绘制该数据集的散点图矩阵和相关系数矩阵图。
- (4) 绘制每加仑油行驶的英里数 (`mpg`)、总马力 (`hp`) 和汽车自重 (`wt`) 3个变量的3D散点图和气泡图。

## 习题2（描述性统计）

本题使用 `exercise4_1.csv`（请在 `rdata` 或 `pydata` 文件夹中寻找）

随机抽取50个网络购物的消费者，调查他们某月的网购金额（单位：元）。

- (1) 计算平均数、标准差、极差和四分位差。
- (2) 计算10%、50%、75%的分位数。
- (3) 计算偏度系数和峰度系数，分析网购金额的分布特征。
- (4) 计算标准分数和极值标准化值，检测数据的离群点。

## 习题3 (选做)

你还知道哪些数据可视化图表，它们在Stata里如何实现？

收集并整理代码，解释图形含义、优缺点、适用场景与案例等。

下面的网站可以帮助你在种类繁多的图表中选择合适的一种：

- 图之典 ↗ (提供matplotlib和ggplot2代码)
- Chart.Guide - helps you choose and design the Perfect Chart ↗
- Data Viz Project | Collection of data visualizations to get inspired and find the right type ↗
- The Data Visualisation Catalogue ↗

## 作业要求

- 提交代码+对应图形+必要解释，注意标明题号
- 整理成一个pdf文档，命名为学号-姓名如“2410000-张三”
- **DDL**：10月23日24点前发送至助教飞书（2120253538）
- 视代码完整性、作图美观性、答案准确性评分
- 逾期提交、抄袭雷同记0分
- 可以使用AI

# 欢迎交流 ~



zzynankai@outlook.com



Bilibili: 西山yu



xishanyu2.github.io