

# R 与统计图形

刘思喆、谢益辉

第三届中国 R 语言会议 © 上海财经大学

2010 年 11 月 14 日

# 目录

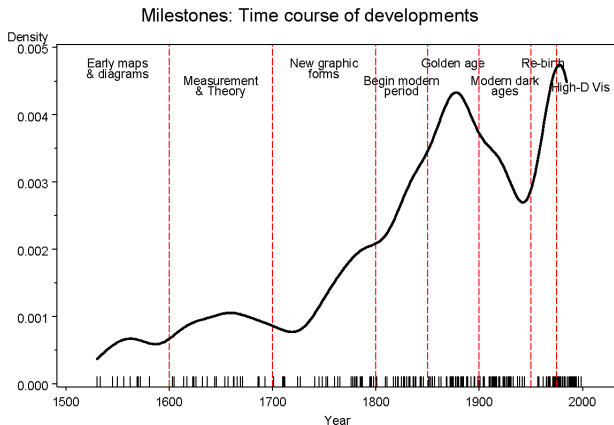
- ① 历史
- ② 输出和设备
  - 输出
  - 设备
- ③ 一般性统计绘图
  - 高级
  - 低级
- ④ 扩展
  - 参数和颜色
  - 通用
  - 特殊
  - 交互
  - 动画
  - 图库
- ⑤ 模型

# 1 历史

# 图形的各个历史时期:

- Pre - 17th Century: Early maps and diagrams
- 1600 - 1699: Measurement and theory
- 1700 - 1799: New graphic forms
- 1800 - 1850: Beginnings of modern graphics
- 1850 - 1900: The Golden Age of statistical graphics
- 1900 - 1950: The modern dark ages
- 1950 - 1975: Re-birth of data visualization
- 1975 - present: High-D, interactive and dynamic data visualization

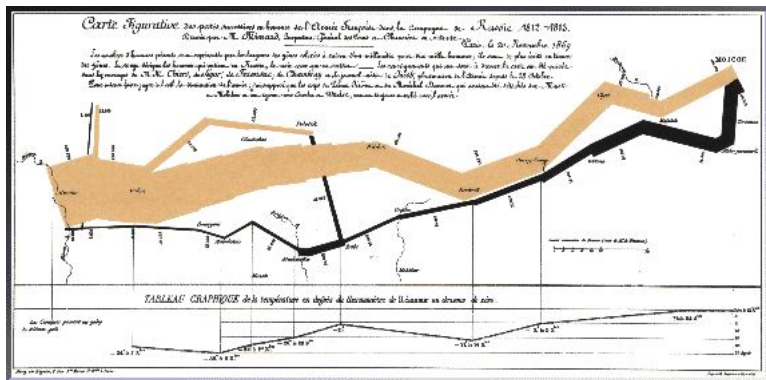
19 世纪后期是数据可视化的黄金时期；20 世纪中叶以后，随着计算机技术的发展，可视化技术又一次重生



数据可视化历史中，各个时期里程碑事件的分布，由密度估计和坐标须展示。

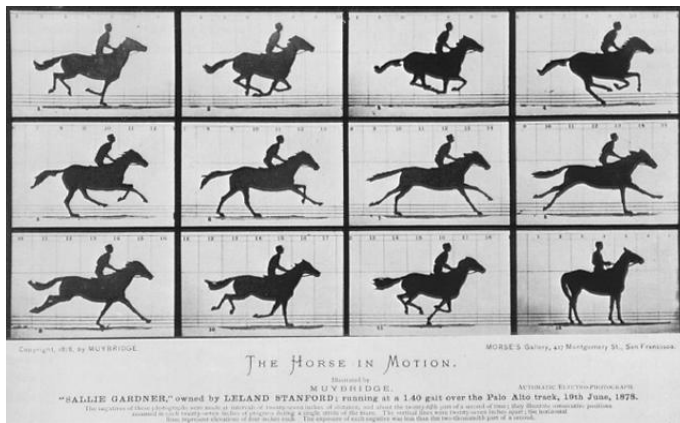
来源: Handbook of Computational Statistics: Data Visualization(2006)

# 绘图历史上最好的图形 (1869): 拿破仑远征俄国



The French engineer, Charles Minardi (1781-1870), illustrated graphically the disastrous campaign of Napoleon against Russia in 1812. The width of the course is proportional to the number of surviving soldiers in the war campaign. In beige for the way in and in black for the return.

# Edward Muybridge 的赛马动画 (1877)



来源: <http://didyouknow.org/how-a-horse-kicked-off-the-movie-industry/>

## 2 输出和设备

- 输出
- 设备



# 图形输出

- 当用户发出绘图命令后，R 会自动弹出图形窗口
- 图形窗口上，可选择 save 的文件类型
- savePlot

# R 基础图形设备

R 图形设备

	名称	描述
屏幕 显示	x11 windows	X 窗口 Windows 窗口
文件 设备	postscript pdf pictex png jpeg bmp xfig win.metafile <sup>a</sup>	ps 格式文件 pdf 格式文件 供 L <sup>A</sup> T <sub>E</sub> X 使用的文件 png 格式文件 jpeg 格式文件 bmp 格式文件 供 XFIG 使用的图形格式 emf 格式的文件

<sup>a</sup>仅在 Windows 下有效

R 同时支持 Cairo 图形设备，这种图形设备可以支持高质量的矢量绘图 (PDF,PostScript and SVG)，高质量的点阵绘图 (PNG,JPEG,TIFF) 或屏幕输出。

### 3 一般性统计绘图

- 高级
- 低级
- 参数和颜色

# plot 类函数

语法: `plot(x, y, ...)`

[1]	<code>plot.acf*</code>	<code>plot.data.frame*</code>	<code>plot.decomposed.ts*</code>
[4]	<code>plot.default</code>	<code>plot.dendrogram*</code>	<code>plot.density</code>
[7]	<code>plot.ecdf</code>	<code>plot.factor*</code>	<code>plot.formula*</code>
[10]	<code>plot.hclust*</code>	<code>plot.histogram*</code>	<code>plot.HoltWinters*</code>
[13]	<code>plot.isoreg*</code>	<code>plot.lm</code>	<code>plot.medpolish*</code>
[16]	<code>plot.mlm</code>	<code>plot.ppr*</code>	<code>plot.prcomp*</code>
[19]	<code>plot.princomp*</code>	<code>plot.profile.nls*</code>	<code>plot.spec</code>
[22]	<code>plot.spec.phase</code>	<code>plot.stepfun</code>	<code>plot.TukeyHSD</code>
[25]	<code>plot.stl*</code>	<code>plot.table*</code>	<code>plot.ts</code>

# 其他统计绘图

## 统计绘图函数

*boxplot*, *bxp*, *cdplot*,  
*contours*, *filled.contour*,  
*fourfoldplot*, *hist*, *image*,  
*matplot*, *mosaicplot*, *persp*,  
*spineplot*, *symbols*, *stripchart*, ...

## 参数

*add* = *TRUE*

*axes* = *TRUE*

# 元素

- 点 (points)
- 线 (lines, abline, segments, arrows)
- 多边形 (rect, polygon, box)

# 辅助

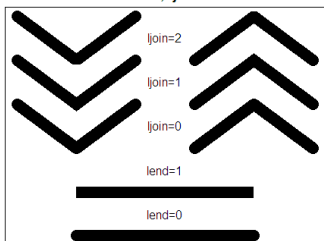
- 坐标轴 (axis)
- 文本 (title, legend, text, mtext)
- 公式及符号 (expression)

# par

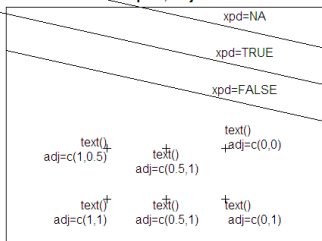
- par 是魔鬼 (72 个参数, 常用的也达到近 20 个)
- 是绘图的必要条件 (除非你想忍受 R 的默认绘图状态)



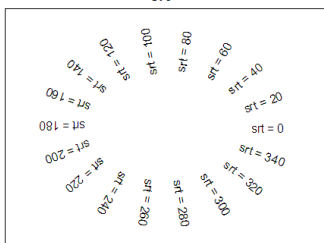
**lend=, ljoin=**



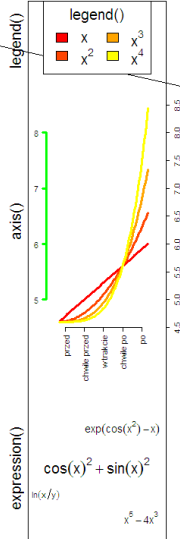
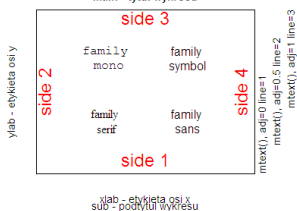
**xpd=, adj=**



**srt=**



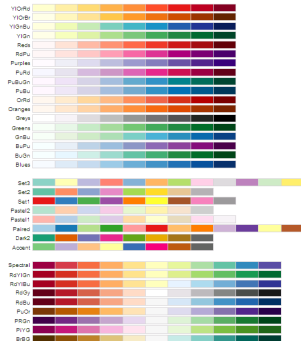
**main, line=**  
**main - tytuł wykresu**



# 颜色

- `colors()`
- 颜色参数的传递
  - 颜色名称, 如 `red`
  - 整数, 对应当前调色板 `palette()`
  - 16 进制的三原色, 如 `#FF0000`
- 特定主题调色板, 如 `heat.colors()`
- 颜色的扩展包, 如 `RColorBrewer`、`colorRamps`

Color names for colors from the RColorBrewer package



Color names for grDevices and colorRamps packages



Color names for colors(grDevices)



# 绘图的例子

plot\_confidence\_interval.r

- 4 扩展
  - 通用
  - 特殊
  - 交互
  - 动画
  - 图库

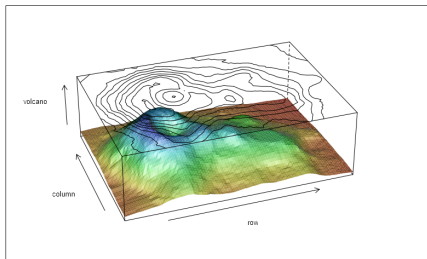
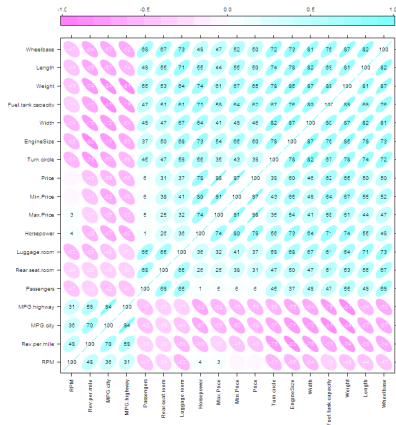
# 通用类

- lattice
- ggplot2

# lattice

- The Trellis graphics of S(+) and The Elements of Graphing Data (Cleveland, 1985).
- The lattice package is an independent implementation of Trellis graphics
- lattice uses Paul Murrell's grid package, which provides more flexible low-level tools

## lattice 绘图示例

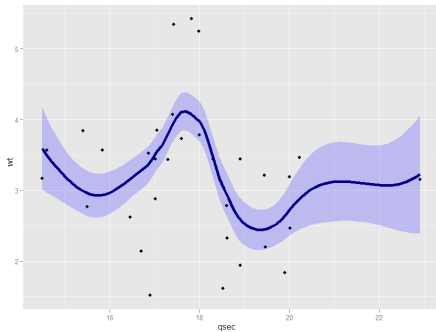
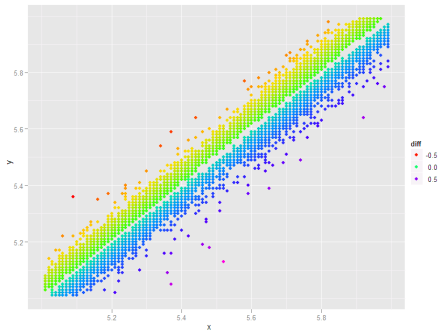




# ggplot2

- The Grammar of Graphics[Wilkinson,2005] 在 R 中的实现
- 图层叠加的概念，如同魔方
  - 几何单位 (Geom, 点? 线? 光滑?)+ 统计变换 (Stat, 直方图? QQ 图?)+ 尺度表示 (Scale, 颜色渐变? 元素大小?)+ 坐标系 (Co-ord, 笛卡尔? 极坐标?)+ 面板分类 (Facet, 根据分类变量分别画图)+ 元素位置调整 (Position, 条形图并列或堆积? 散点随机微小打乱?)
  - 扩展了泛型函数: + (使用非常形象)
- 细节设置自动化，例如图例

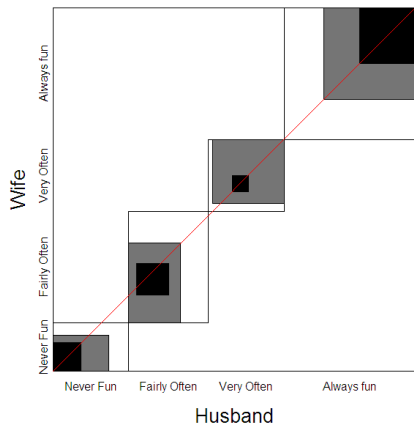
# ggplot2 绘图示例



## vcd

- Visualizing Categorical Data
- The package was inspired by the book "Visualizing Categorical Data" by Michael Friendly.

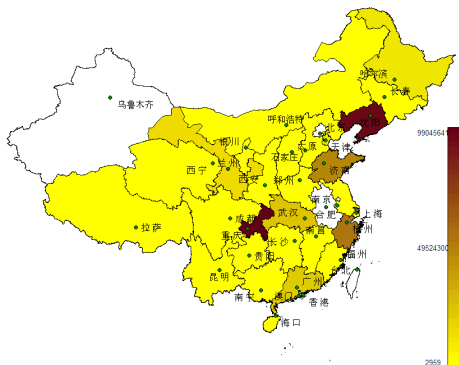
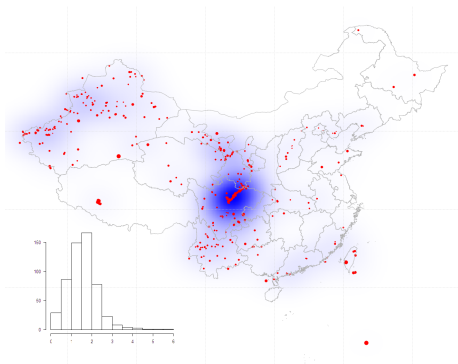
# Bangdiwala's Observer Agreement Chart



	Never Fun	Fairly Often	Very Often	Always fun
Never Fun	7	7	2	3
Fairly Often	2	8	3	7
Very Often	1	5	4	9
Always fun	2	8	9	14

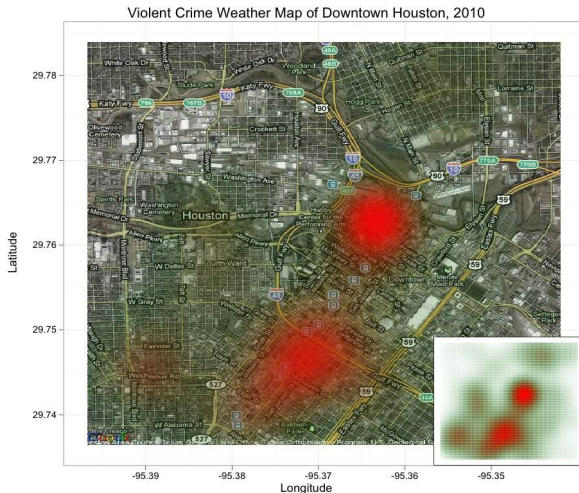
Data from Hout et al. (1987) given by Agresti (1990) summarizing the responses of married couples to the questionnaire item: Sex is fun for me and my partner: (a) never or occasionally, (b) fairly often, (c) very often, (d) almost always.

## maps 系列



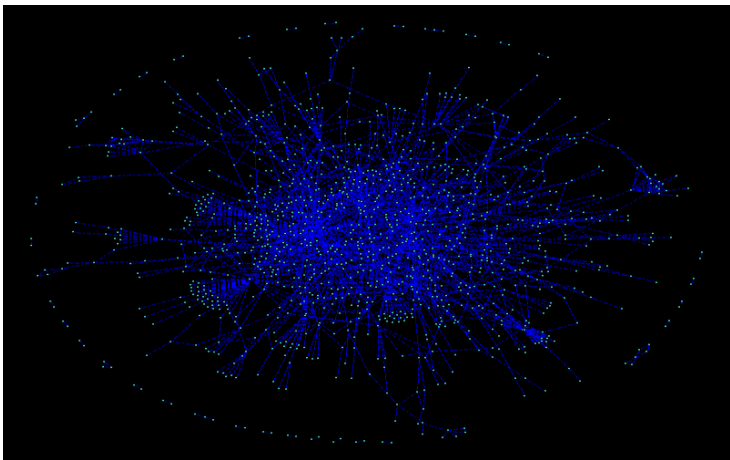
2010 年 11 月 5 日至 10 日国内地震情况以及某企业全国市场份额可视化。其他请参考：  
<http://www.bjt.name/2010/01/chinese-earthquake-visualization/>

# Google Map



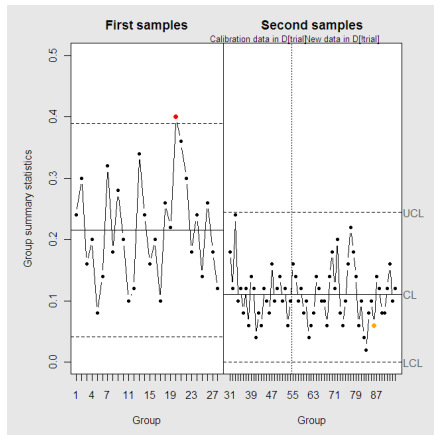
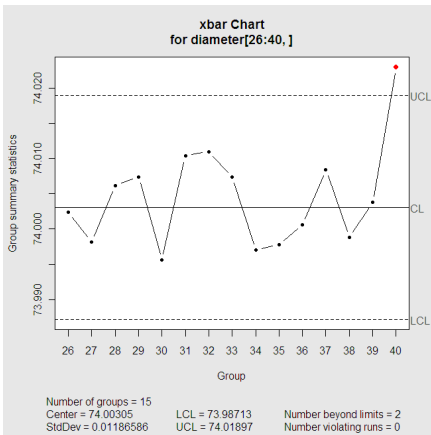
来源: <https://github.com/hadley/ggplot2/wiki/>

## igraph



CRAN 上 R 包的网路关系: <http://www.bjt.name/2009/09/package-networks/>

# Quality Control Charts





# 其他

高维数据展示相关:

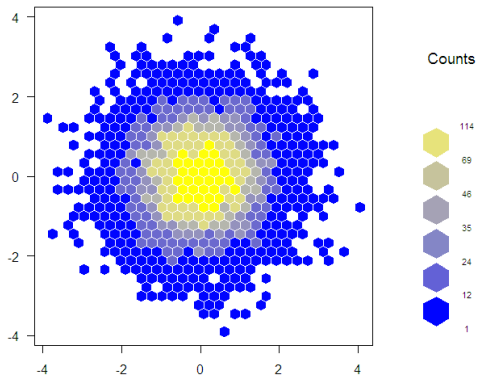
几何投影 (Geometric projection)

轮廓图, Andrews 调和曲线, ...

基于图标 (Icon-based)

脸谱图, 星图, 雷达图, 花瓣图, ...

大数据集



# 交互式绘图

- playwith
- rgl
- rggobi
- iplots
- 例子

# 动画相关的包

如果你有足够的想象力，可以自己来设计适合自己的动画！  
或者借助一些包：

- `animation`(推荐)
- `write.gif` in package `caTools`
- `EBImage` package in BioConductor

# 各类基础图形

- 一维数据：条形图、饼图、Cleveland 点图、坐标轴须、带状图
- 散点图：散点图、向日葵散点图
- 曲线：函数曲线
- 密度和分布：直方图、茎叶图、QQ 图
- 汇总：箱线图、因素效应图
- 分类数据关联：关联图、马赛克图
- 分类对连续：条件密度图、棘状图
- 分类画图：协同图
- 三维图形：颜色图、等高线图、三维透视图、平滑散点图
- 高维散点图：散点图矩阵、符号图

## R 的图库

<http://addictedtor.free.fr/graphiques/>

Read my blog / RELATED SITES: R-project / CRAN / Bioconductor / R-Movies gallery /

search :

Home Browse Related Source code Graphics List Thumbnails

a agreement analysis and association back bar barplot boxplot boxplots chart classification cluster color colored colors conditional conditional conditioning contour explor correlation cumulative curve curves d dem dendrogram density diagram distribution double ellipses escape estimator extended filled fonts for function geographic hershey heibin highest histogram in kernel lattice map mathematical matrices matrix maunga model more mosaic of parallel perspective pie plot plots plotting quiver T regions regression rgb roc rose sample scatter scatterplot seasonal sequences simple spin space special splom treaplot ternary the tree us use validation vector violin volcano volcanoes walter wave whau wheel wind winner wireframe with

Enter the gallery

- 3 best ranked graphs
- 3 last included graphs
- 3 random selected graphs
- 3 last seen graphs

**R Project**

R is a system for statistical computation and graphics. It consists of a language plus a run-time environment with graphics, a debugger, access to certain system functions, and the ability to run programs stored in script files.

**R Graphic Engine**

One of R's strengths is the ease with which well-designed publication-quality plots can be produced, including mathematical symbols and formulae where needed. Great care has been taken over the defaults for the minor design choices in graphics, but the user retains full control.

**Thanks**

Thanks to the R Core Team for the wonderful work made on our favorite tool.

**Technology**

This site is coded in XHTML with CSS. It uses MySQL and PHP. [This Site](#) has been used for the design.

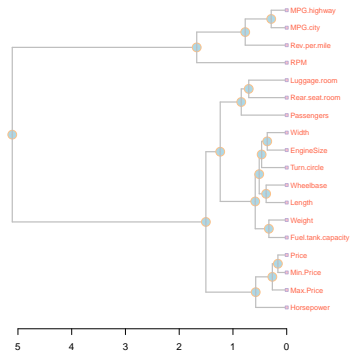
**Related Websites**

- R-project
- CRAN
- Bioconductor
- r-spatial
- R Techniques and Graphics Gallery
- Paul Huet's book
- Michael Friendly's Gallery of Data's Best and Worst of Statistical Graphs

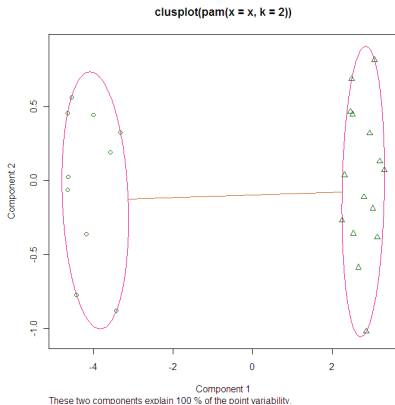
## 5 模型

## 聚类

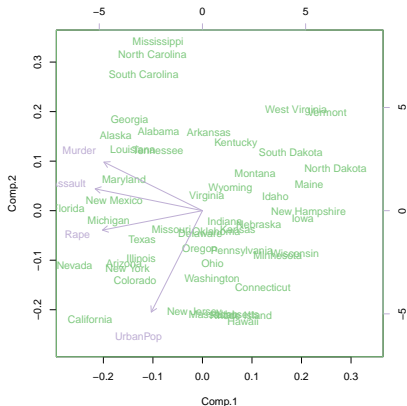
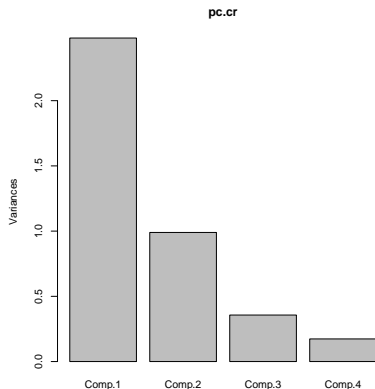
## Hierarchical Clustering



## Partitioning Around Medoids



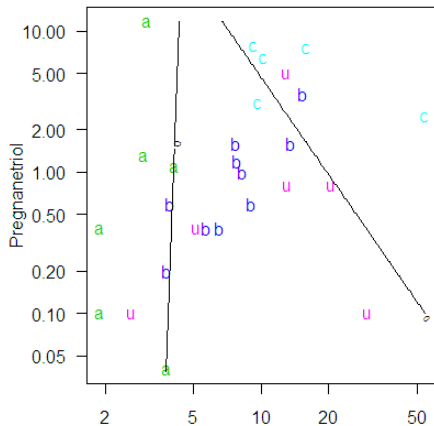
# 主成份和因子分析



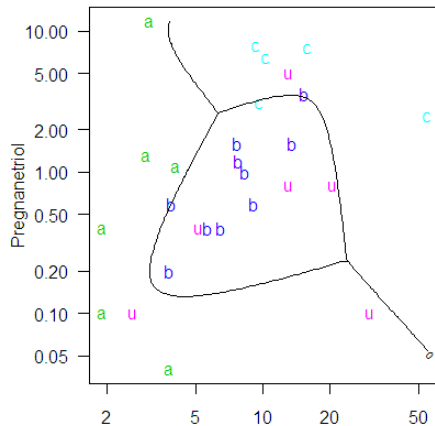


## 判别分析

LDA

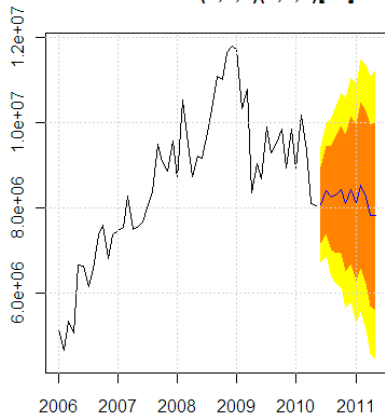


QDA

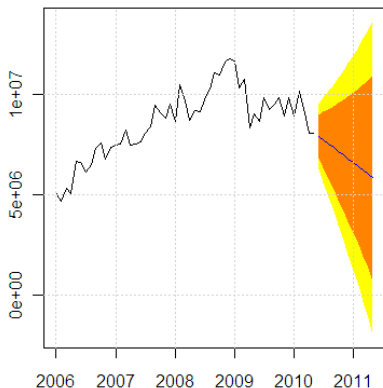


# 时间序列

### Forecasts from ARIMA(1,1,0)(1,0,0)[12]

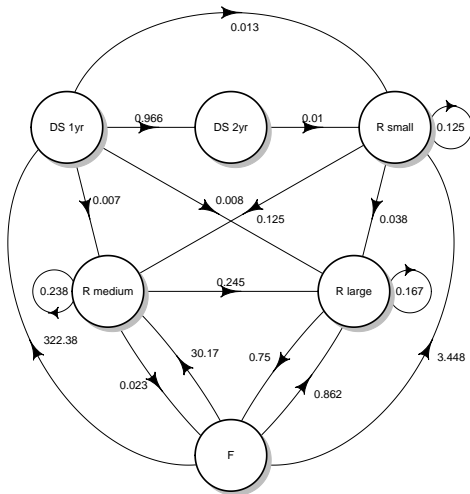


### Forecasts from HoltWinters



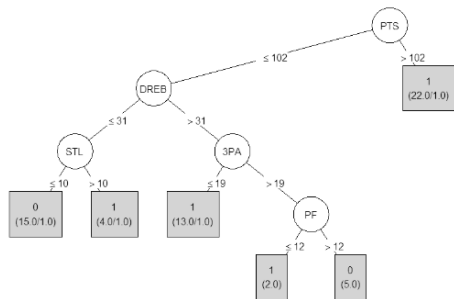
## 图模型

## Life cycle of teasel



# 树型模型

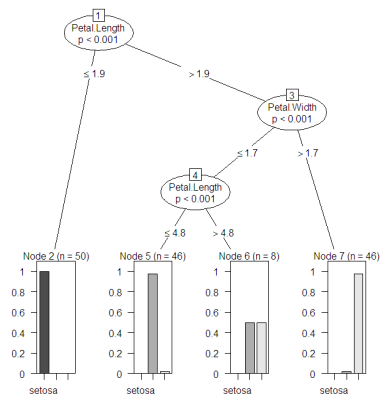
J48



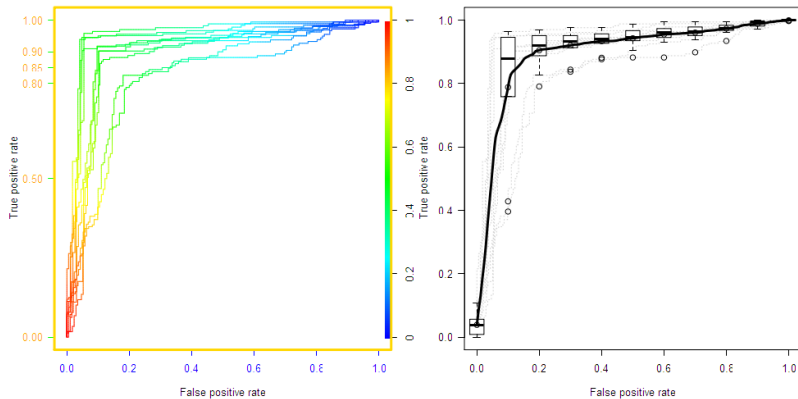
//bjt.name/2009/03/data-2008-2009-rockets

h

## Conditional Inference Trees



# ROC 曲线



# 参考文献



W. N. Venables and B. D. Ripley  
*Modern Applied Statistics with S* 2002



John M. Chambers  
*Software for Data Analysis: Programming with R* 2007



M. Friendly  
*A Brief History of Data Visualization* 2005



Deepayan Sarkar  
*Lattice Multivariate Data Visualization with R* 2007



Leland Wilkinson  
*The Grammar of Graphics*, 2<sup>nd</sup> ed 2005



谢益辉  
现代统计图形 2010

# 关于

## 刘思喆

- Email: [sunbjt<at>gmail.com](mailto:sunbjt@gmail.com)
- Blog : <http://www.bjt.name>

## 谢益辉

- Email: [xie<at>yihui.name](mailto:xie<at>yihui.name)
- Blog : <http://www.yihui.name>

## 问题请关注:

- COS 论坛 R 版: <http://cos.name/cn/forum/15>
- 谢谢各位!

[Jump to first slide](#)