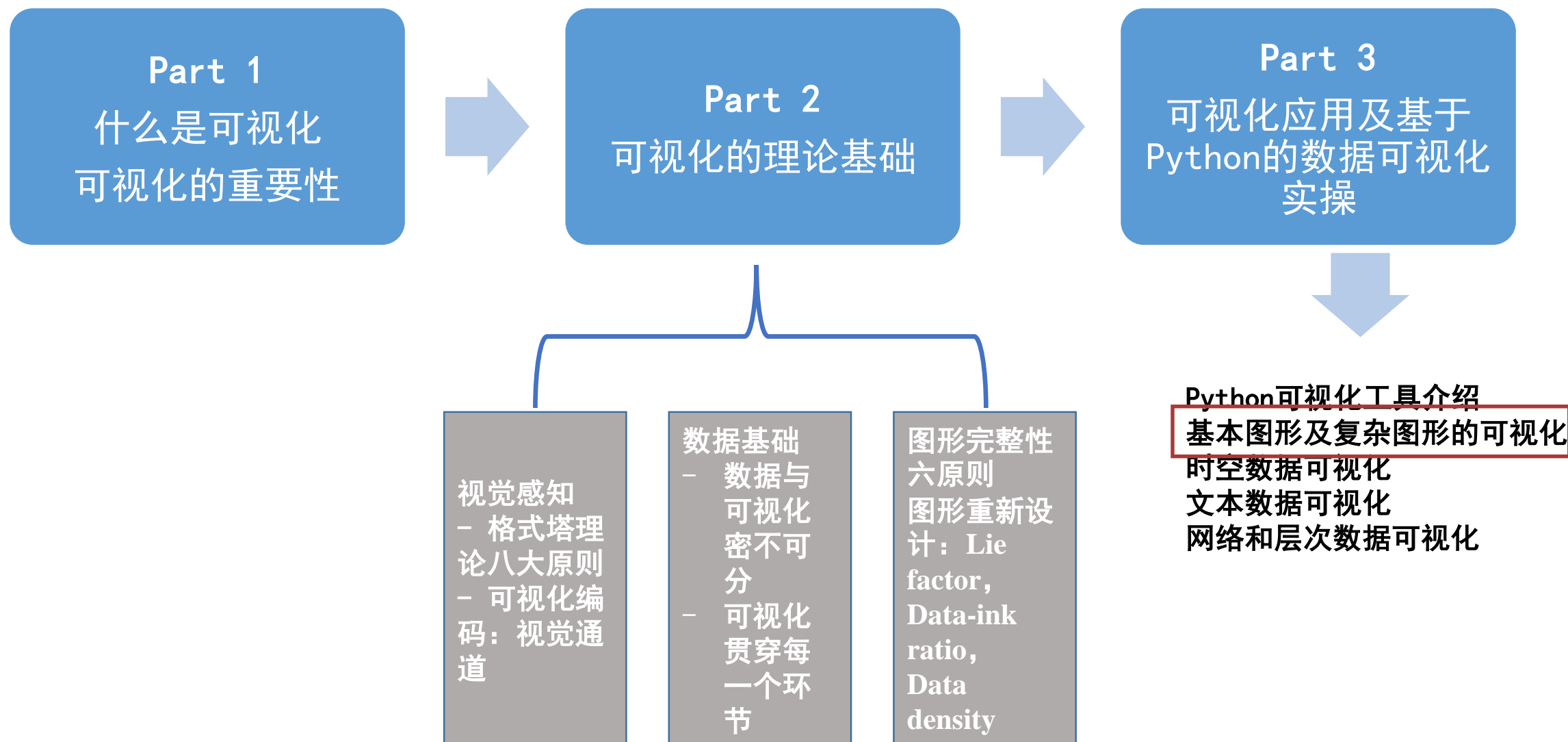




数据可视化

第五讲：可视化方法-基本图形

Our structure



可视化任务

- 迭代的数据探索过程



目录

- 数据变换
- 统计图表

数据变换

- 标准化
- 平滑化
- 采样
- 分箱
- 降维
- 聚类

数据变换

- 目的
 - 更好地解决特定问题
 - 提供更多的可视化设计选择
- 举例 – 数值型温度
 - 发现温度变化的异常值 – 数值型温度
 - 分析全球温度是否升高 – 数值型温度差
 - 判断水温是否适合洗澡 – 序数型(hot, warm, cold)

数据归一化

目的

根据分布映射数据

颜色/尺寸/坐标位置编码

归一化区间：

$[-1, 1]$

$[0, 1]$

数据变换

线性变换

$$y = \frac{x - MinValue}{MaxValue - MinValue}$$

[0, 1]

反正切变换

$$y = \frac{\arctan(x) * 2}{\pi}$$

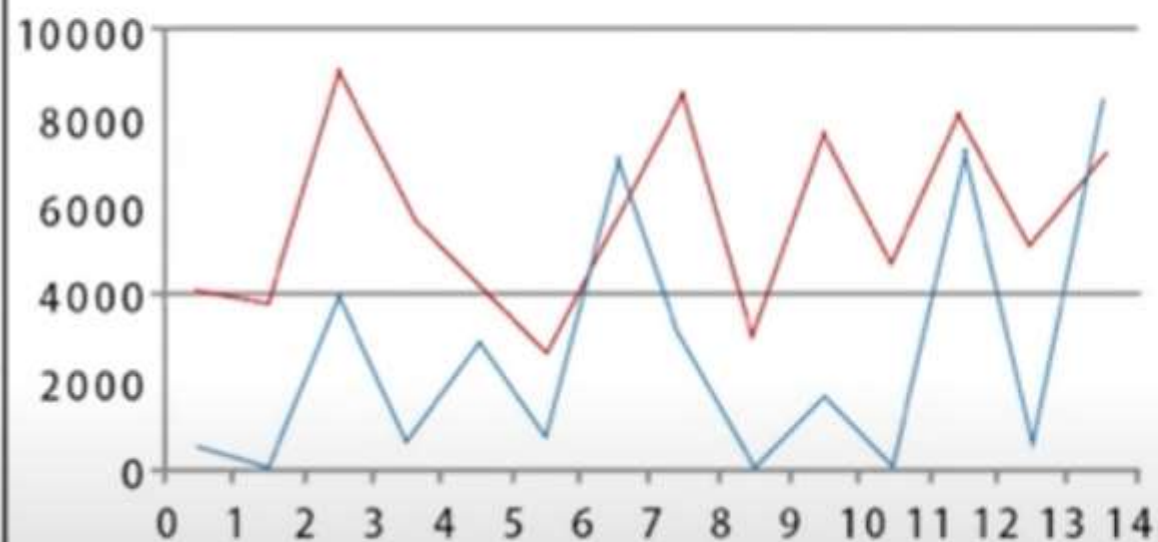
[-1, 1]

对数变换

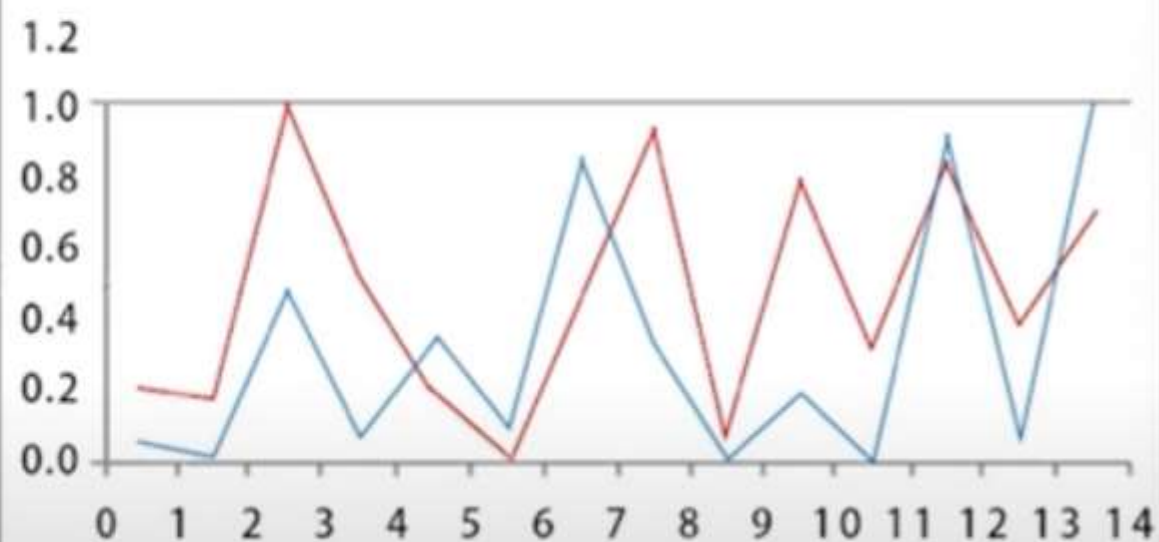
$$y = \log_{10}(x)$$

可以自定义变换函数

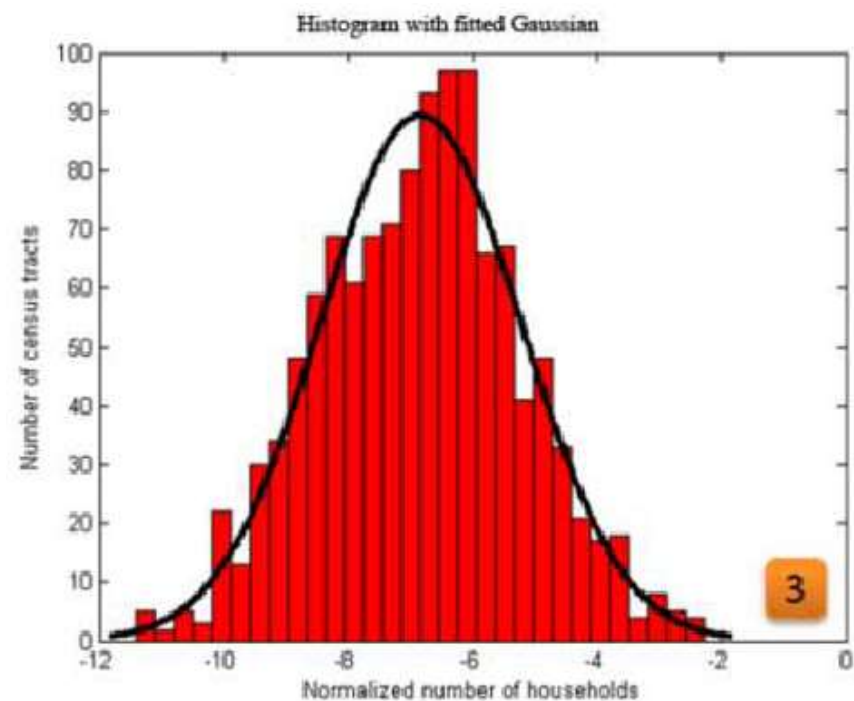
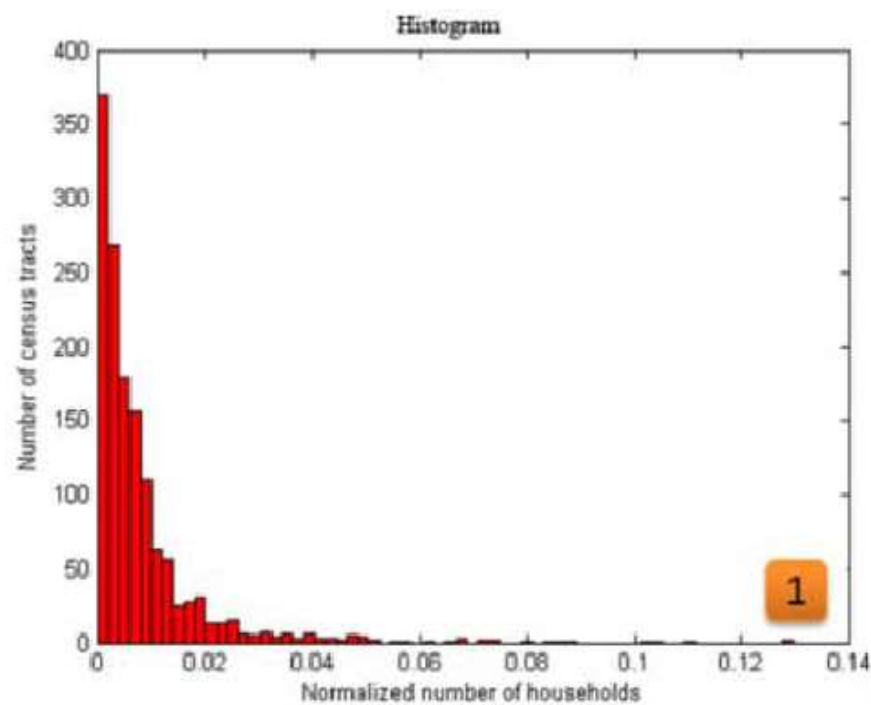
标准化处理前效果图



标准化处理后效果图

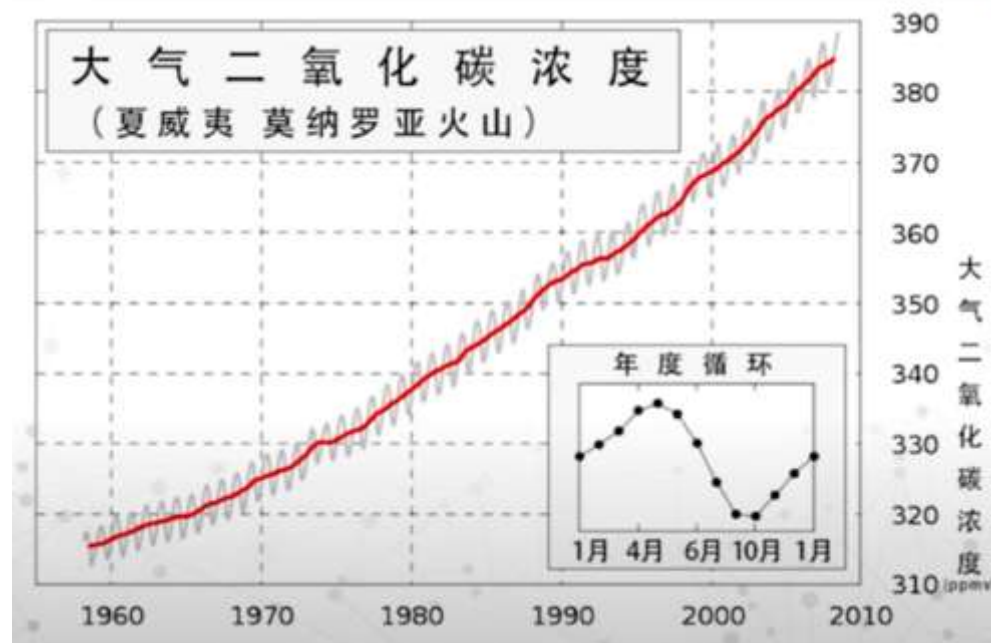
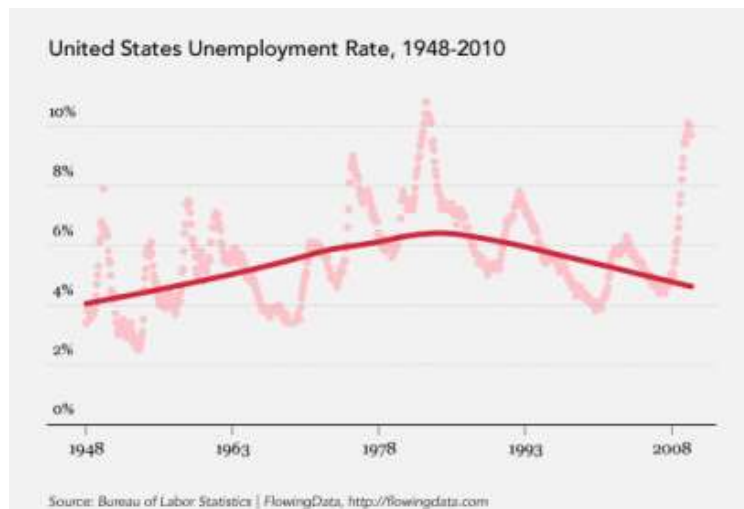


指数变换



曲线拟合/数据光滑化

- 将数据转化成平滑连续的曲线
 - 将数据从“微小的细节”中转移到“更高层面的趋势观察和判断”
- 目的：展示数据趋势



曲线拟合

- 不同的拟合方式：表达并观测趋势（劫富济贫）

- 线性回归

$$\min_{\vec{x}} \sum_{i=1}^n (y_m - y_i)^2.$$

- PLSR (partial least squares regression, 偏最小二乘拟合)
- LOESS (Locally weighted scatterplot smoothing)

- 指数函数

- 多项式曲线

- 自定义方程

统计采样

什么是统计采样？

从统计分布中选出的样本

用于近似原分布中的特征

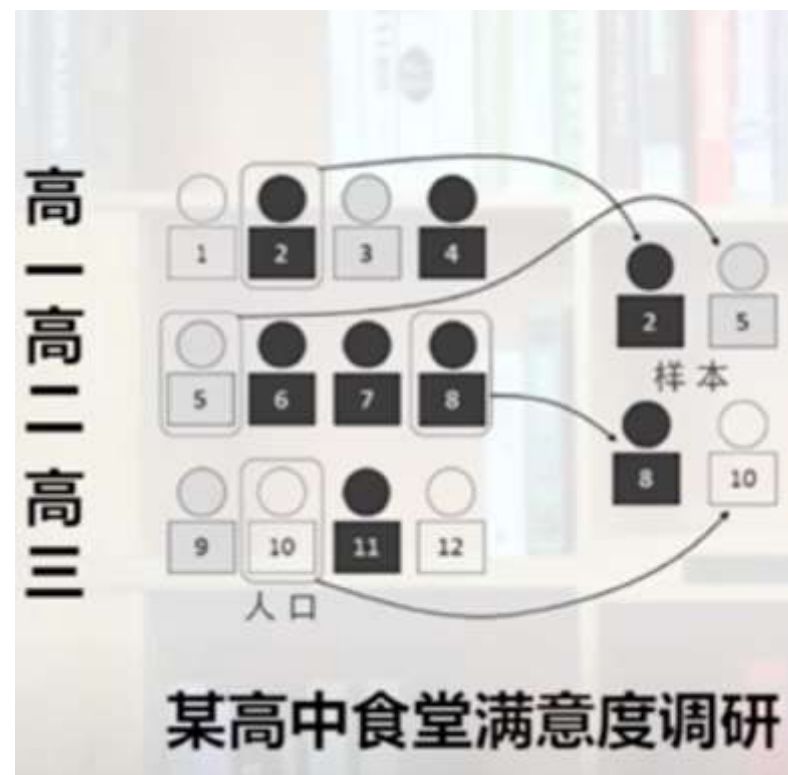
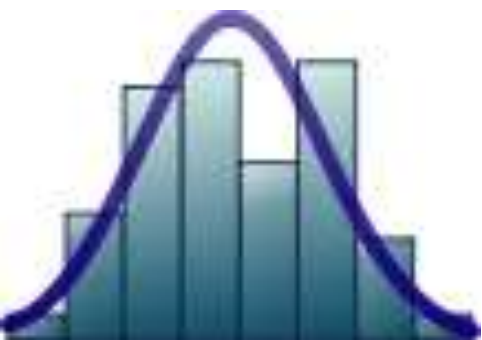
影响采样的因素

分布本身的特性

数据的测量精度

是否需要分析样本细节（样本精细度）

采样成本

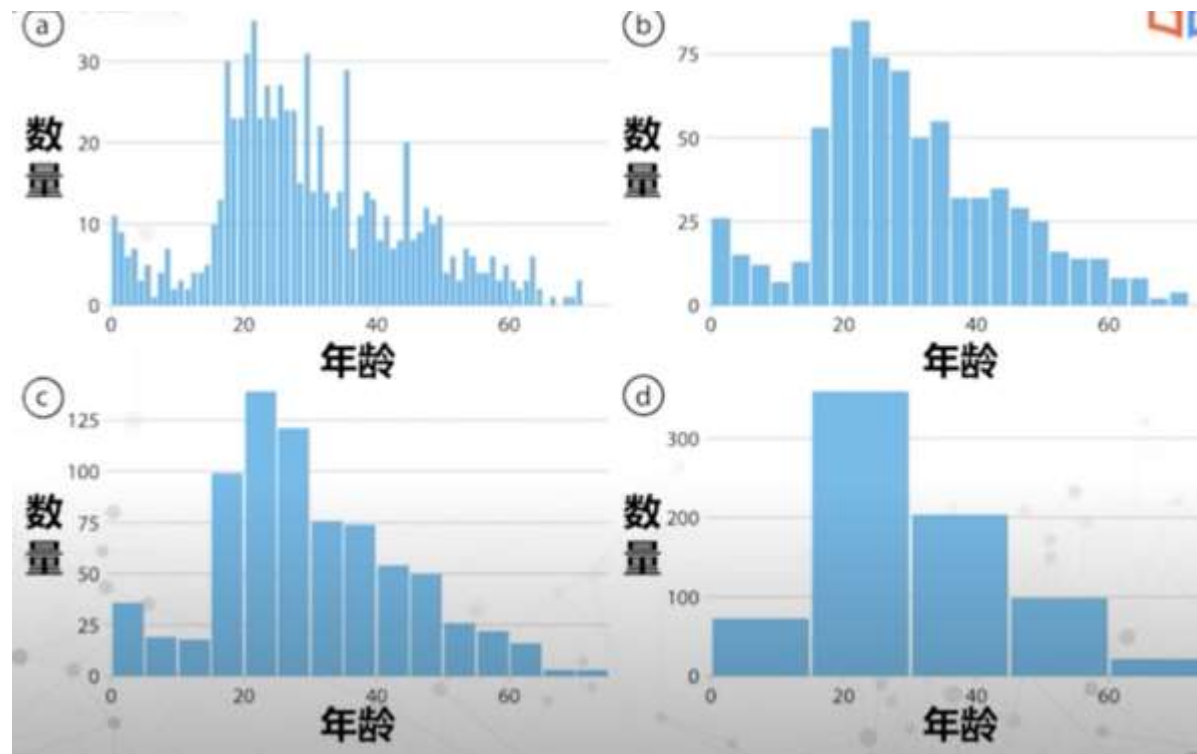


分箱

- 连续数据离散化
 - 将一些连续值分组装进一些“小箱子”的方法
- 单一维度 -> 分箱
- 多个维度



数据降维

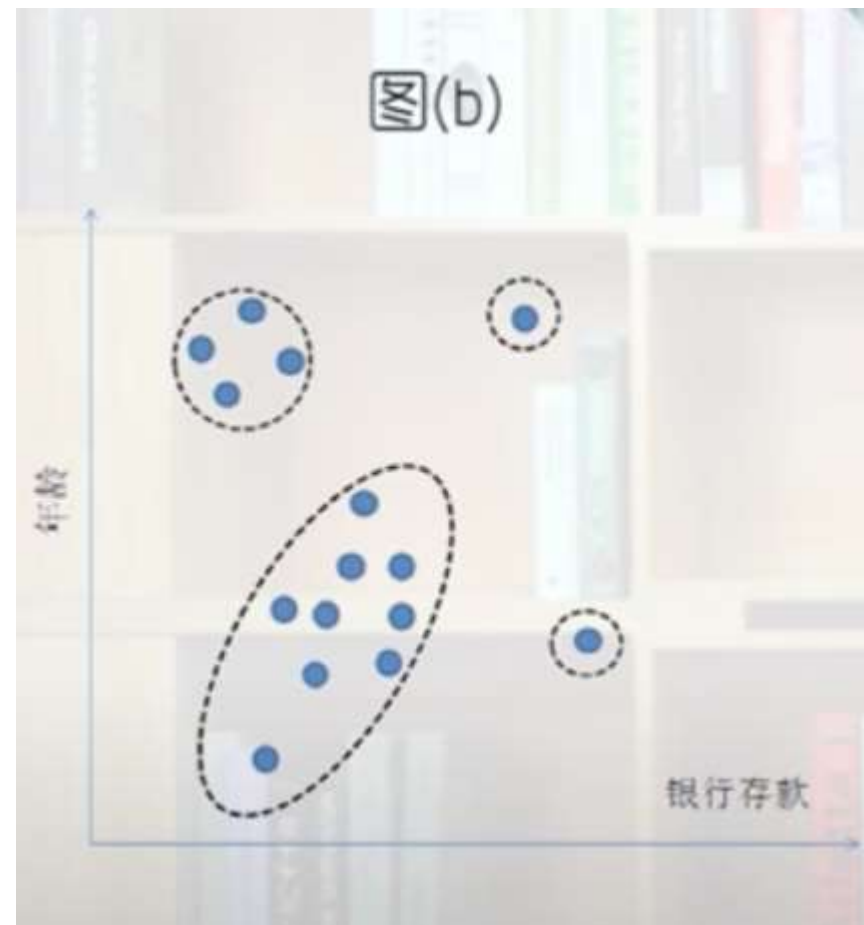
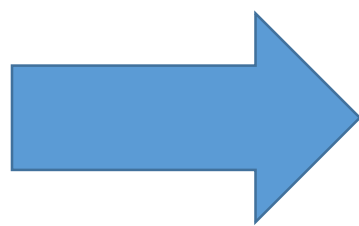
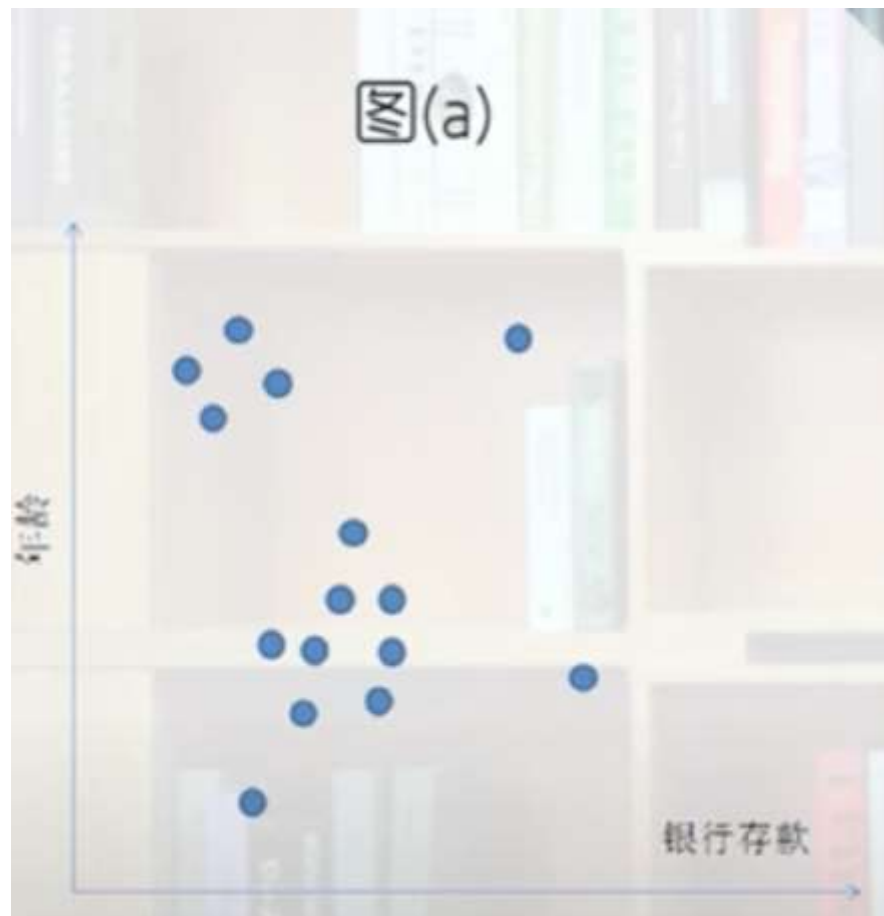


数据降维*

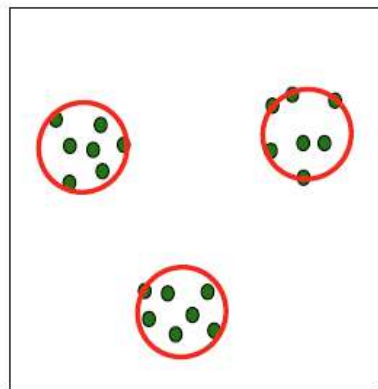
把数据从多维的空间投影到二维或者三维的空间

- 线性方法，如：
 - 主元分析（Principal Components Analysis, PCA）
 - 多维尺度标记（Multidimensional Scaling, MDS）
- 非线性方法，如：
 - T分布随即近邻潜入（t-SNE）
 - 自组织网络（Self-Organizing Map, SOM）
 - 等距特征映射（ISOMAP）

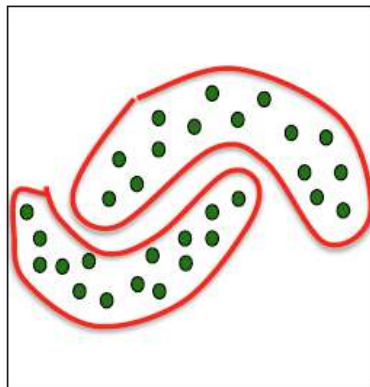
数据聚类



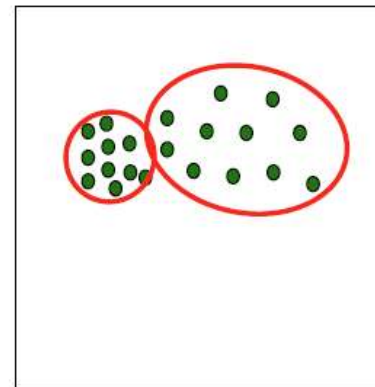
聚类方法



Location



Shape



Density

聚类方法：

K-means聚类

Expectation-Maximization Clustering (EM) *

Gaussian Mixture Model (GMM)*

Spectral Clustering*

Hierarchical Clustering*

K-Means聚类

- K-means

随机产生K个中心位置

将每个数据点归为距离最近的中心位置所属的类

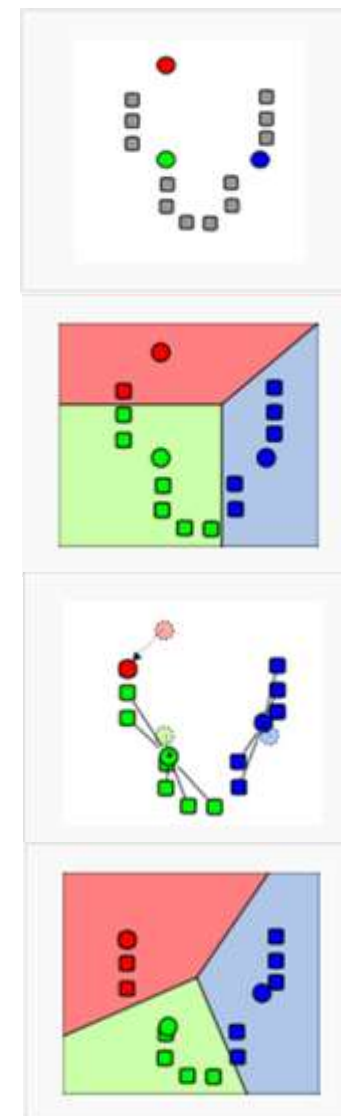
根据新的类别划分重新计算中心位置

回到第二步，直到满足一定约束

- K-medoids – 改进

中心位置必须在数据点所在位置

中心位置满足“到类内所有数据点的距离之和最小”



目录

- 数据变换
 - 标准化
 - 平滑化
 - 采样
 - 分箱
 - 降维
 - 聚类
- 统计图表：原始数据绘图

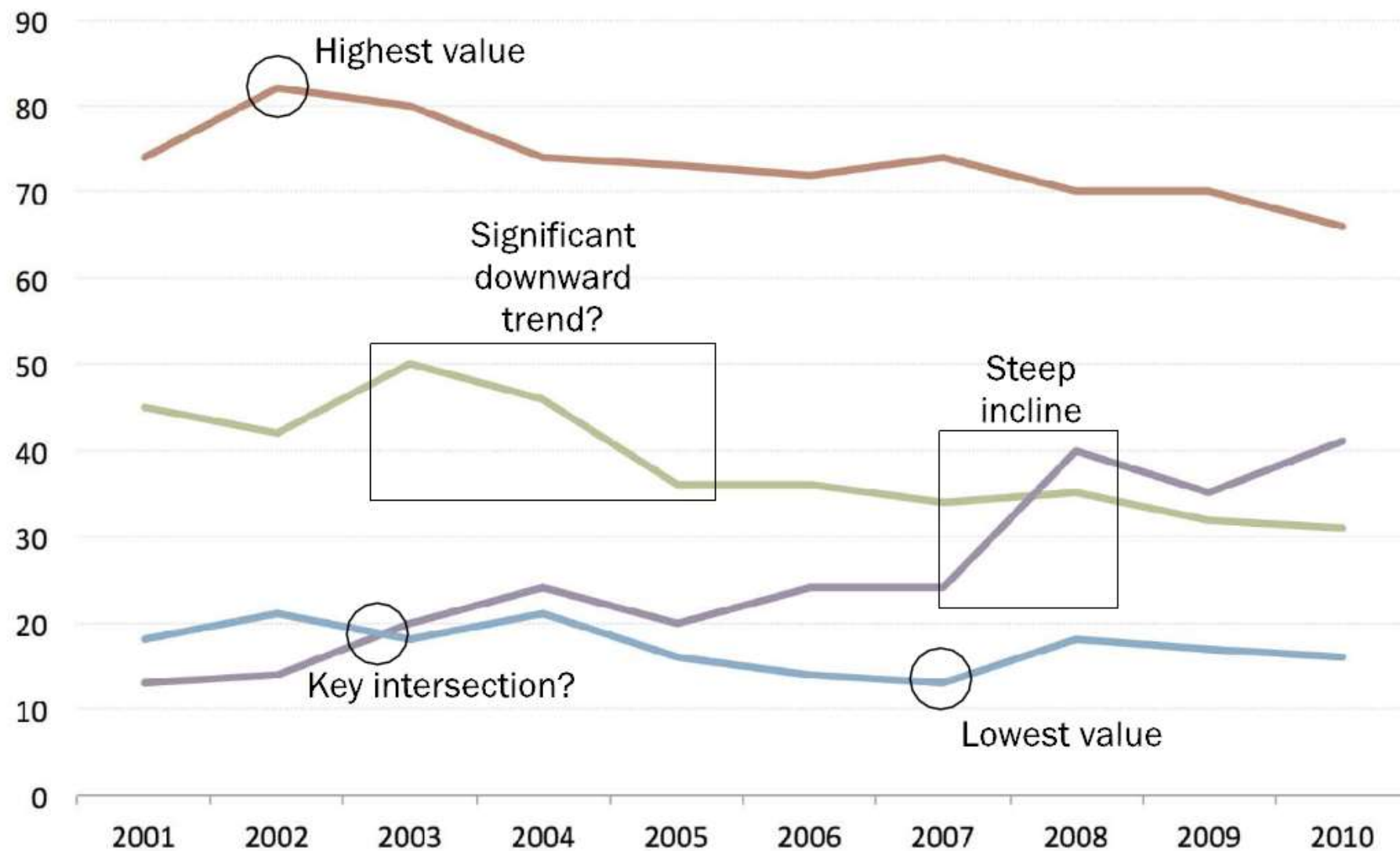
统计图表

- 作用
 - 比较与比例
 - 趋势和模式
 - 关系

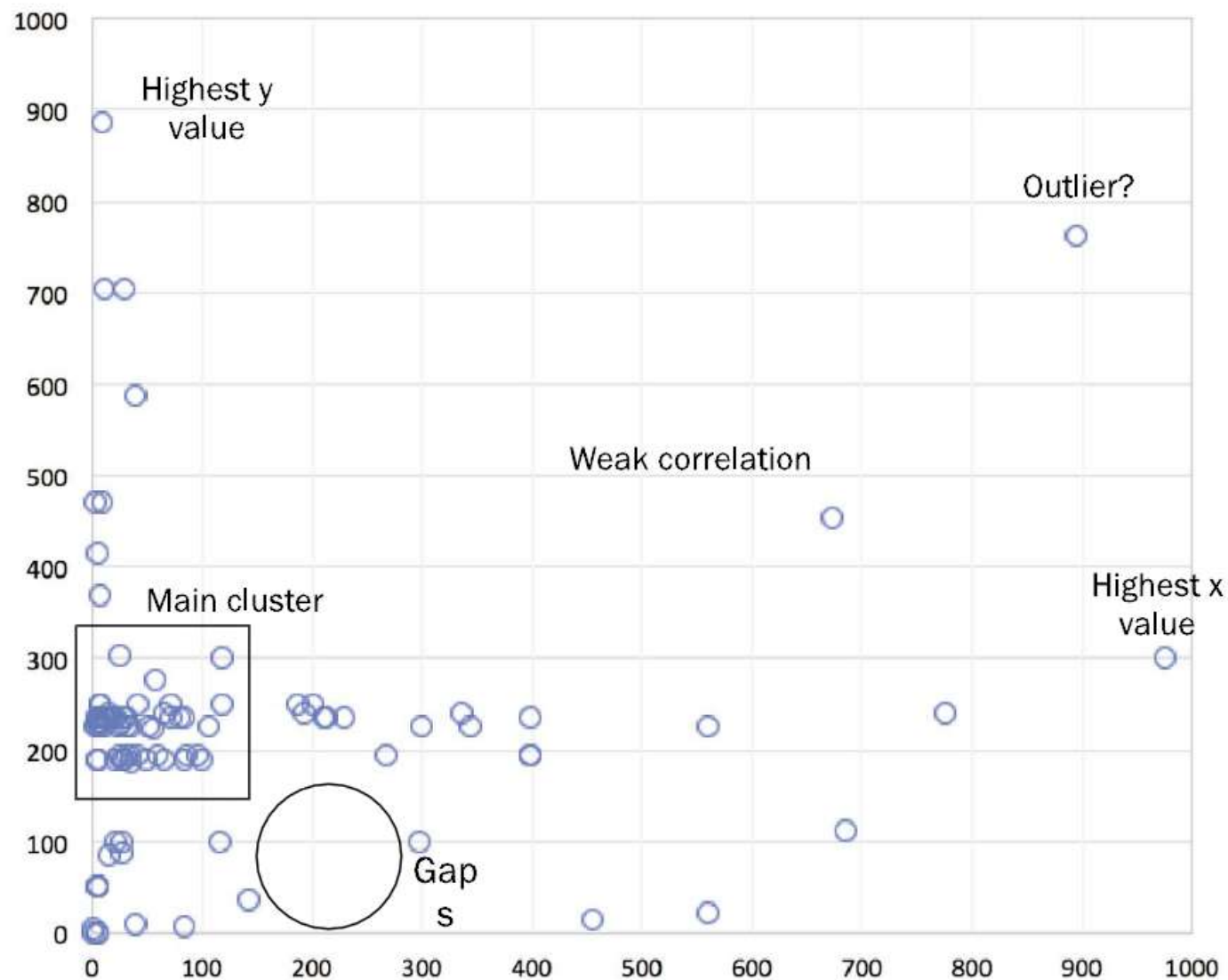
统计图表 – 比较与比例



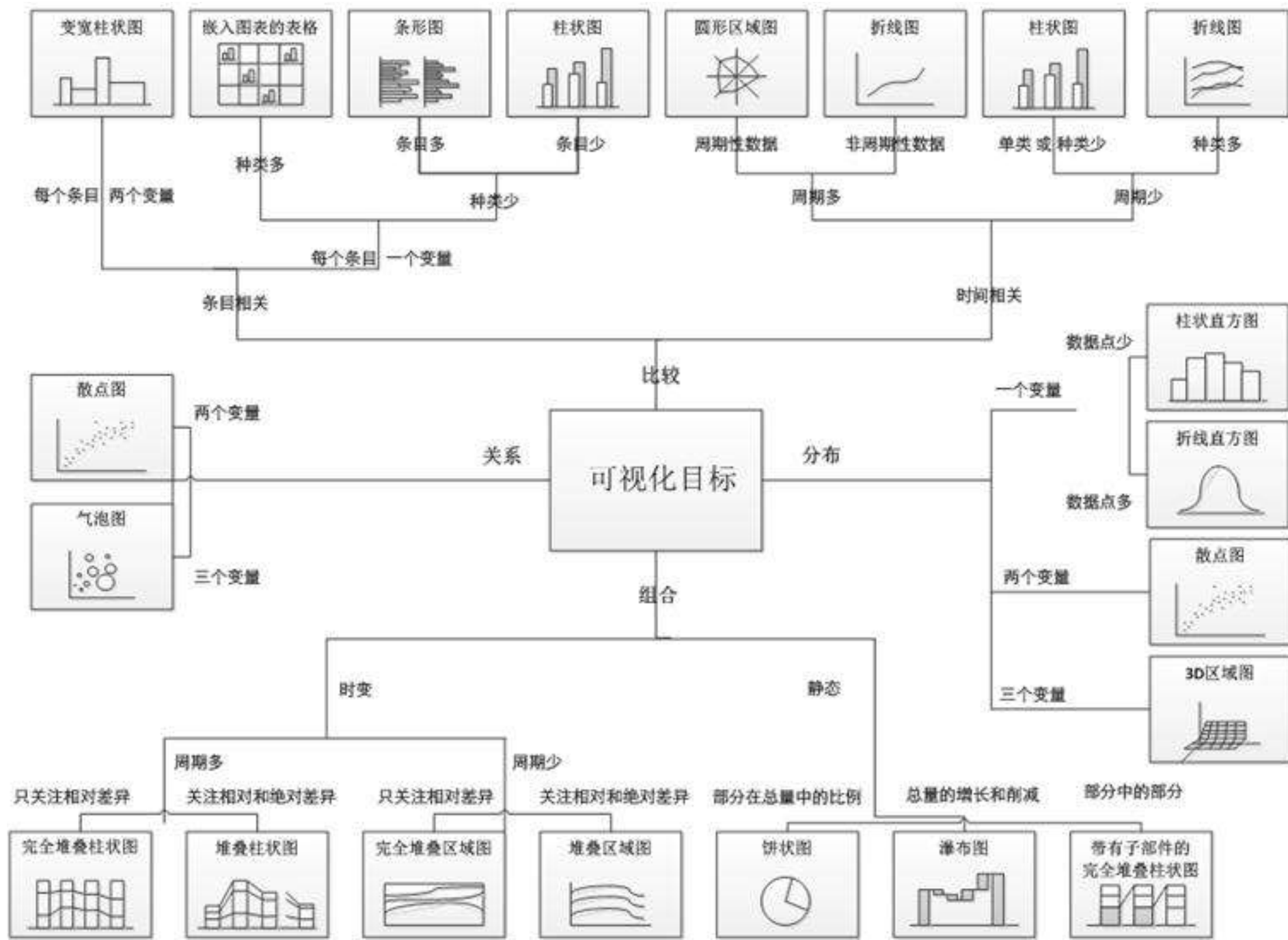
统计图表 - 趋势与模式



统计图表 - 关系

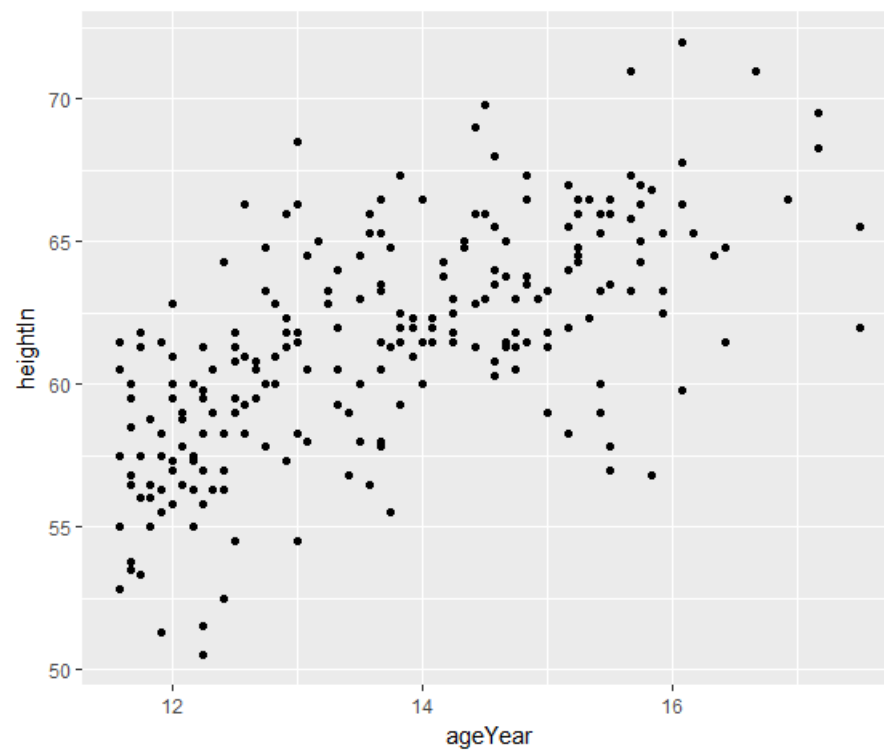


统计图表



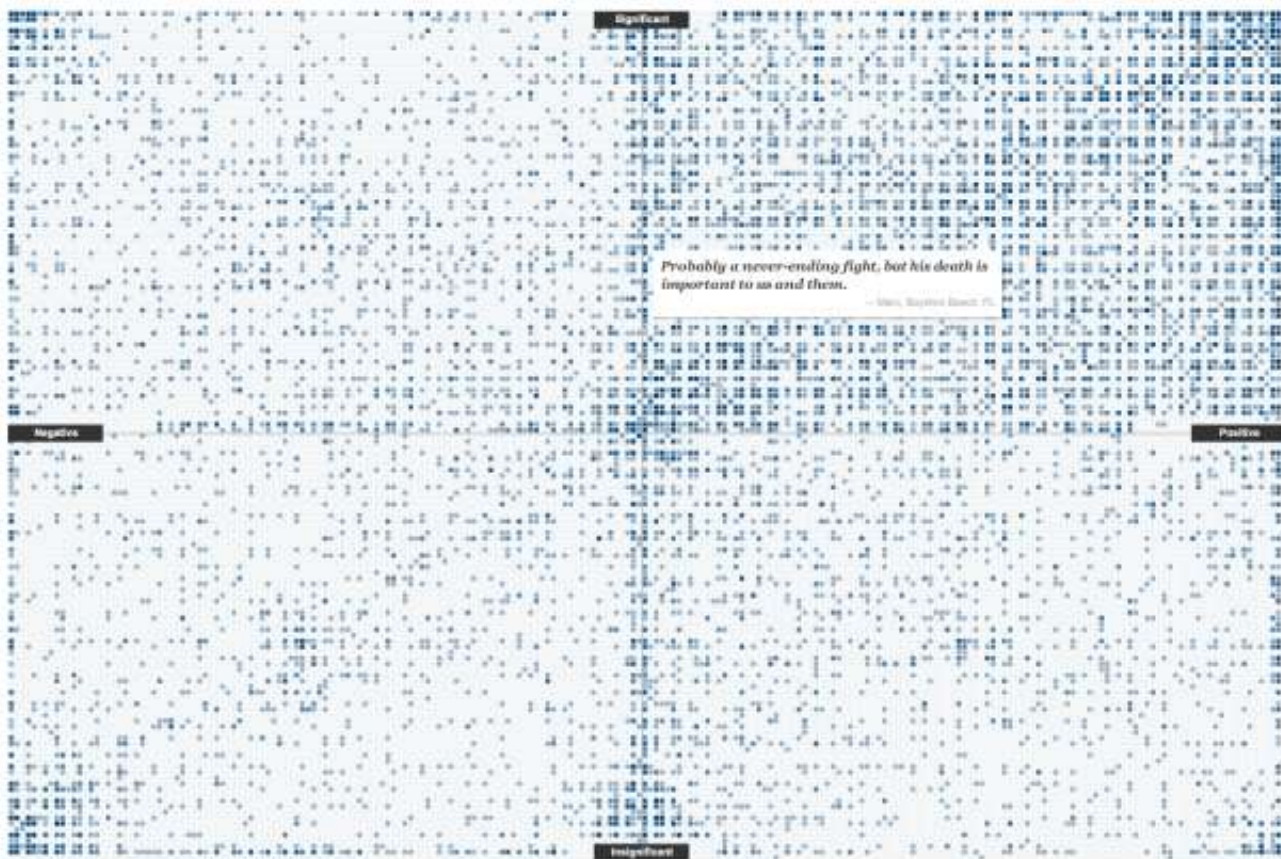
关系：散点图

- 散点图是表示二维数据的标准方法，探讨变量之间关系



本拉登之死

意义重大



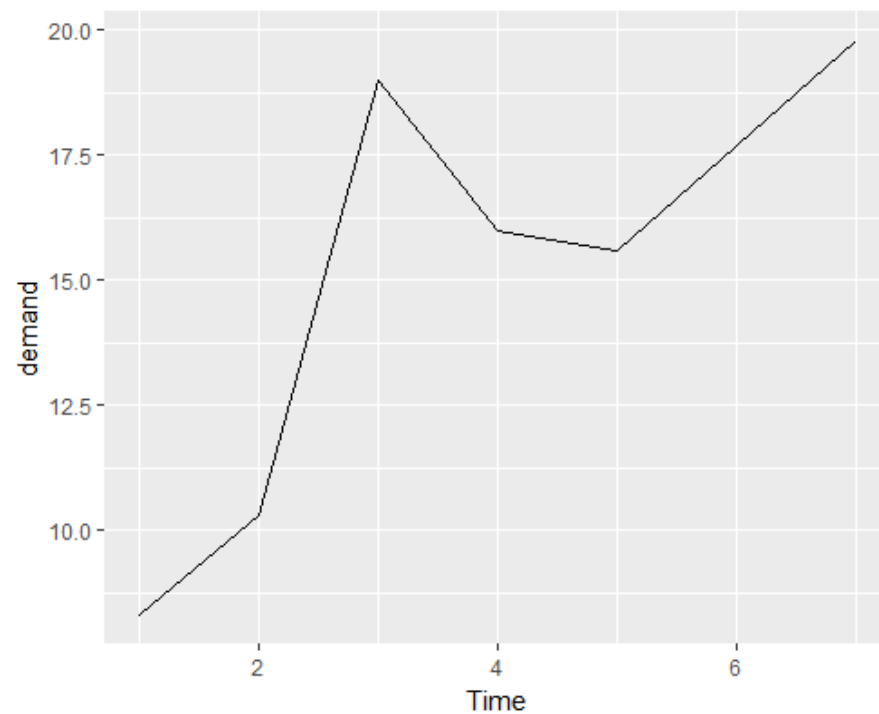
坏事

好事

毫无意义

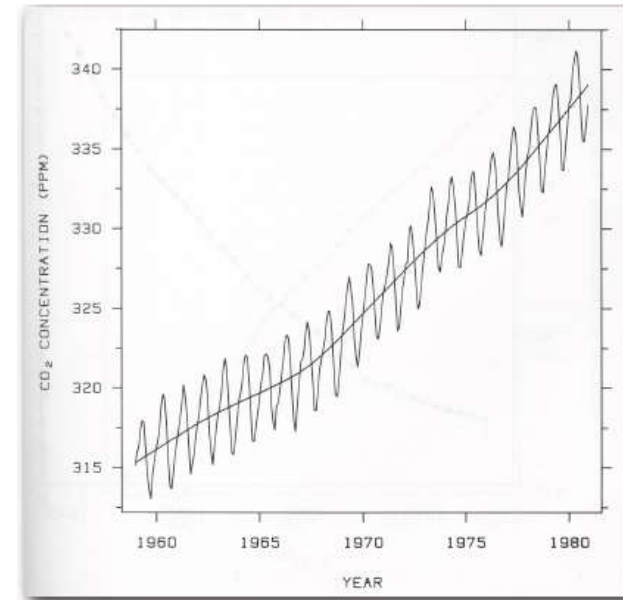
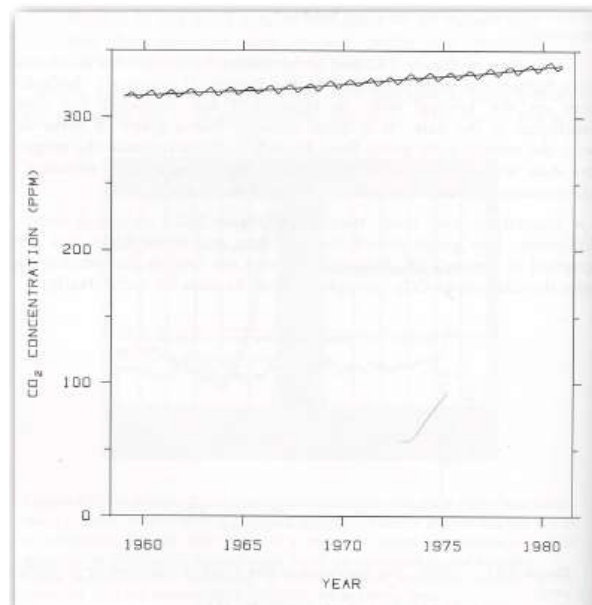
比较：折线图

- 使用直线段来连接一系列点



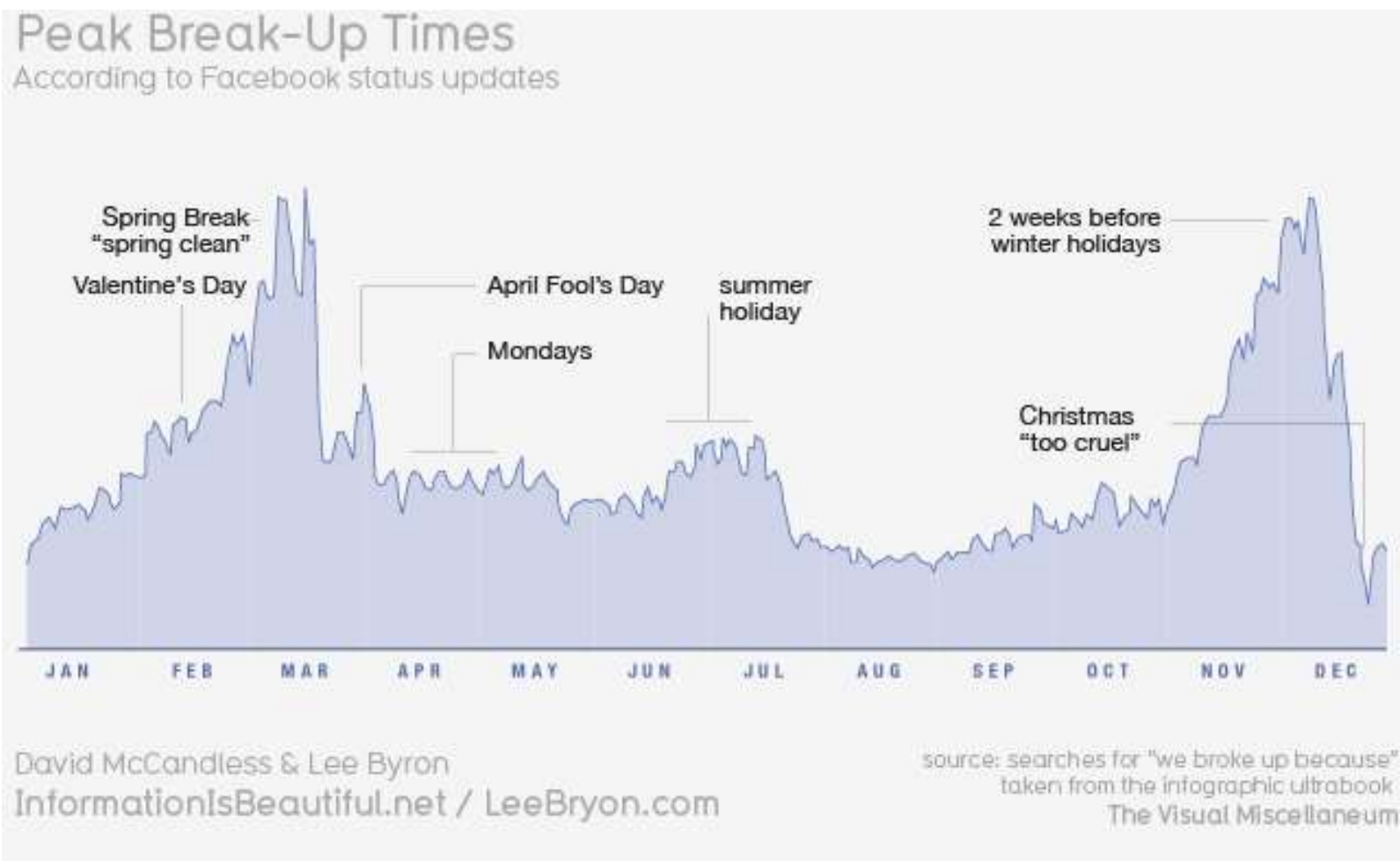
比较：折线图

- 宽高比



Multi-Scale Banking to 45°, 宽高比：左7.87，右1.17

折线图



统计图对比



折线图

重量级
同时表达数据走势和分布

Sparkline

轻量级
只表达数据走势



By Edward Tufte

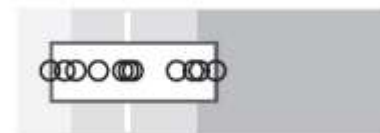
Bandline

中量级?
表达数据走势和分布

数据走势 + 模糊分布



数据分布



By Stephen Few

比较：柱状图/条形图

- 采用长方形的形状和颜色编码数据的属性



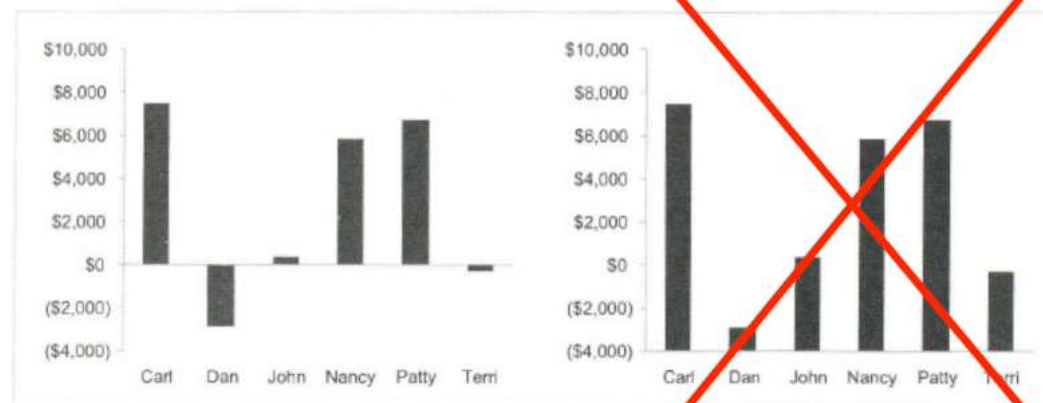
每个国家消费了多少啤酒？

每人每周消耗的瓶数

柱状图

- 注意事项

- 尺度
- 使用零点为基准点
- 不必要的三维设计



堆叠柱状图

- 每根直柱内部也可用像素图方式编码
 - 分解整体，比较局部信息



每个柱子代表一个整体

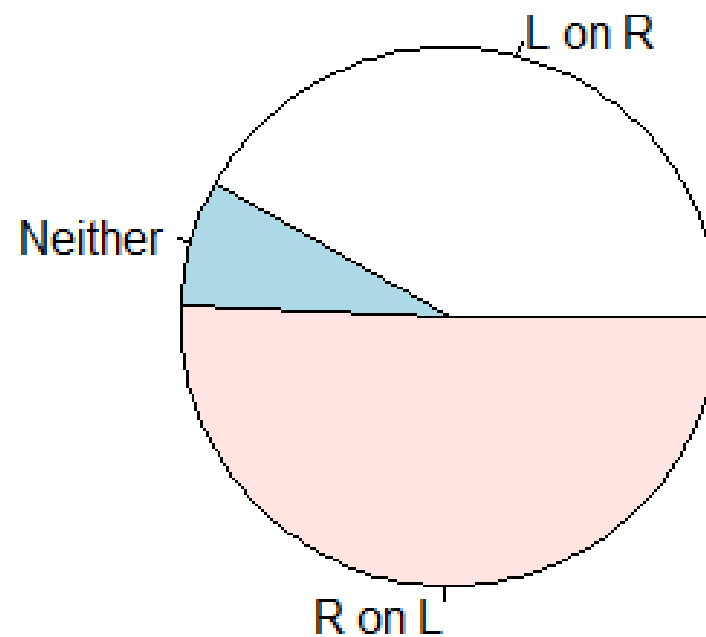
柱子中不同颜色代表不同类别

堆叠柱状图

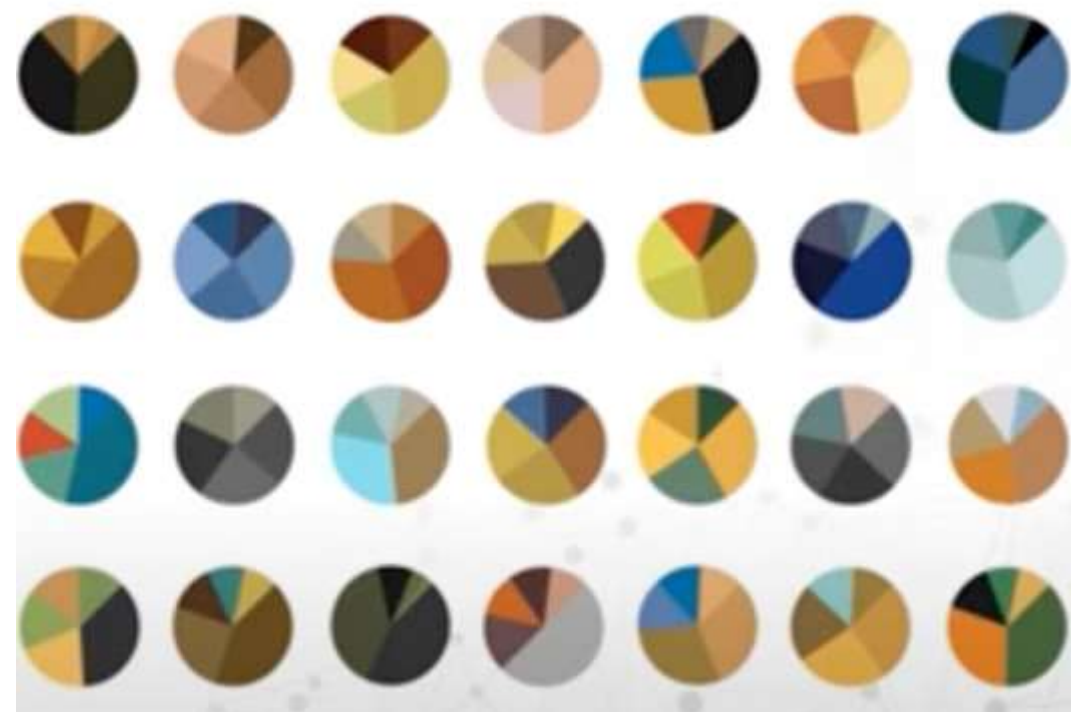
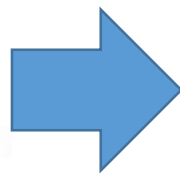


饼图

- 用环状方式呈现各分量在整体中的比例，是环状树图的基础
 - 采用了饼干的隐喻
 - 避免3D

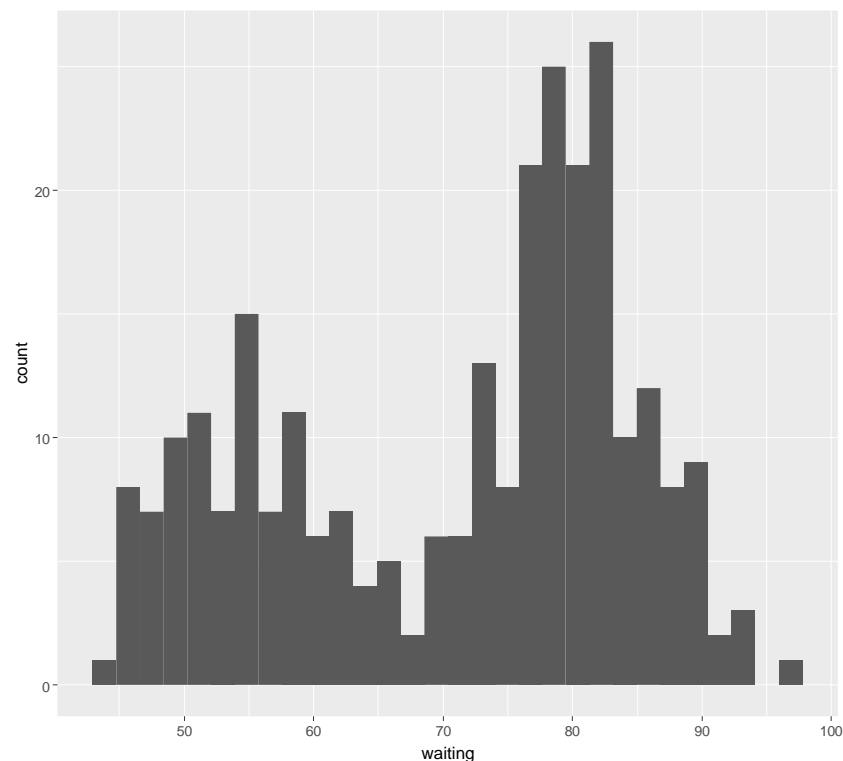


饼图：梵高的作品



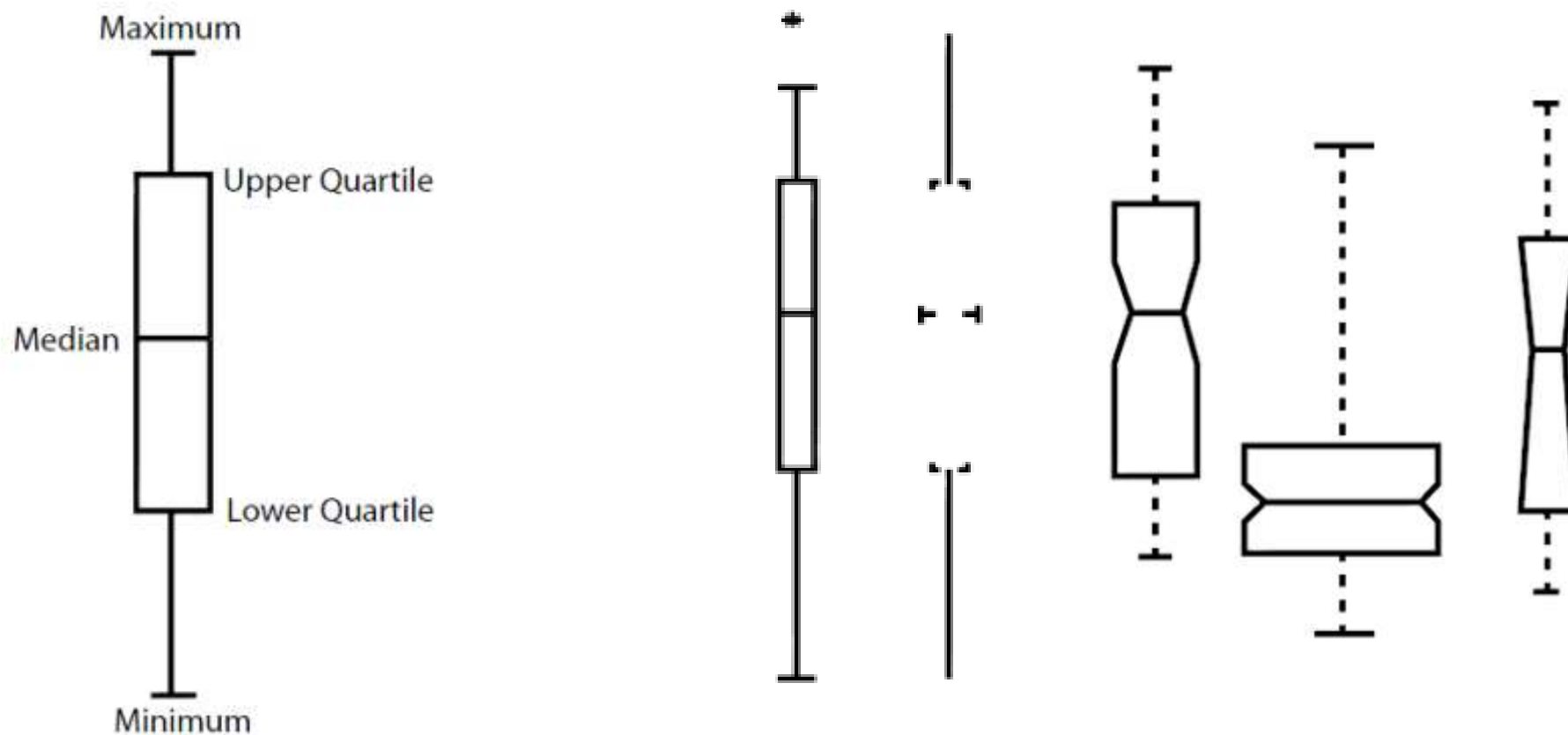
分布：直方图

- 对数据集的某个数据属性的频率统计
- 呈现数据分布、离群值和分布的模态

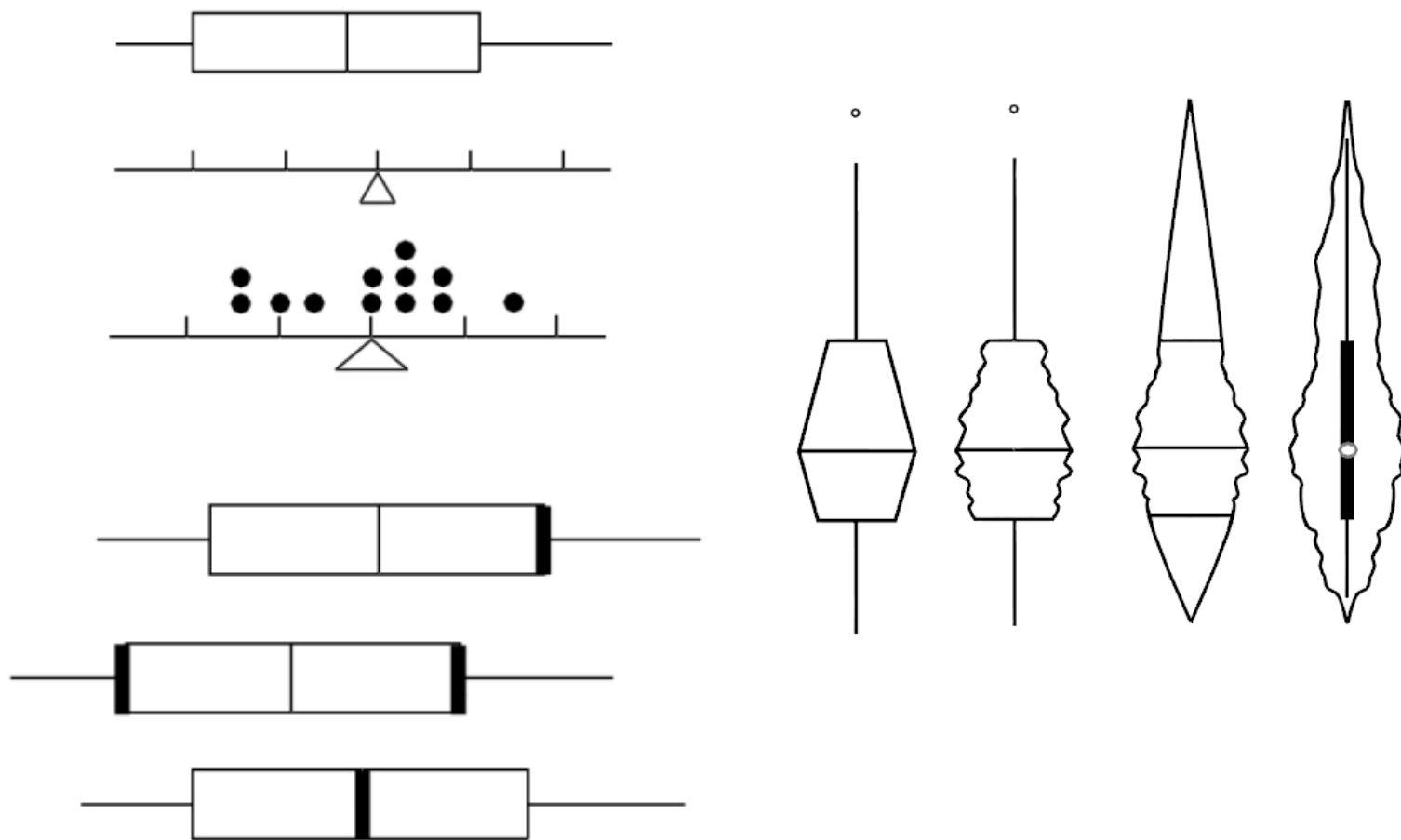


分布：箱线图/盒须图

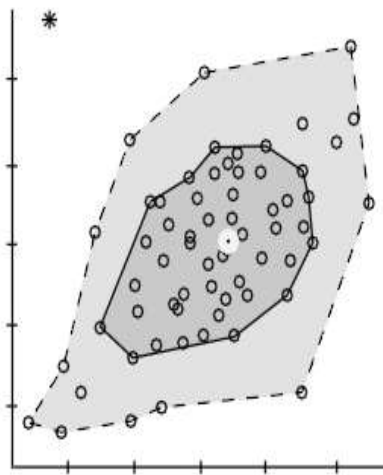
- 也叫盒须图。长方形盒子表示数据的大概范围



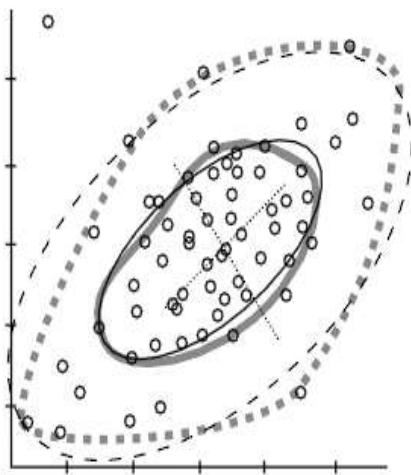
盒须图变种



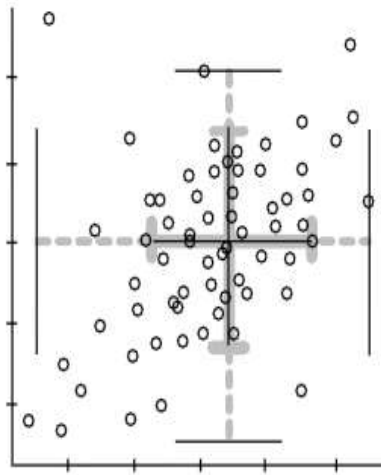
盒须图变种



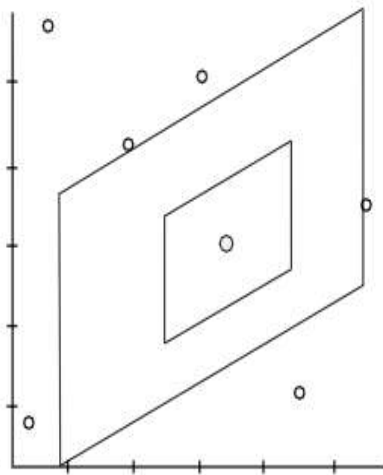
2D Box Plot



Relplot



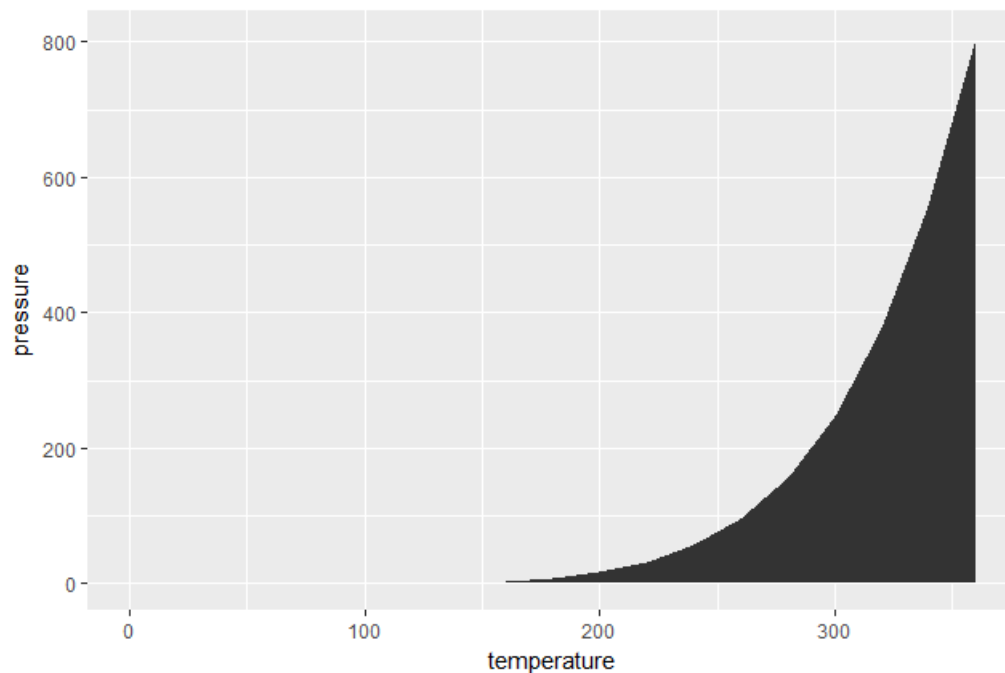
Rangefinder Box Plot



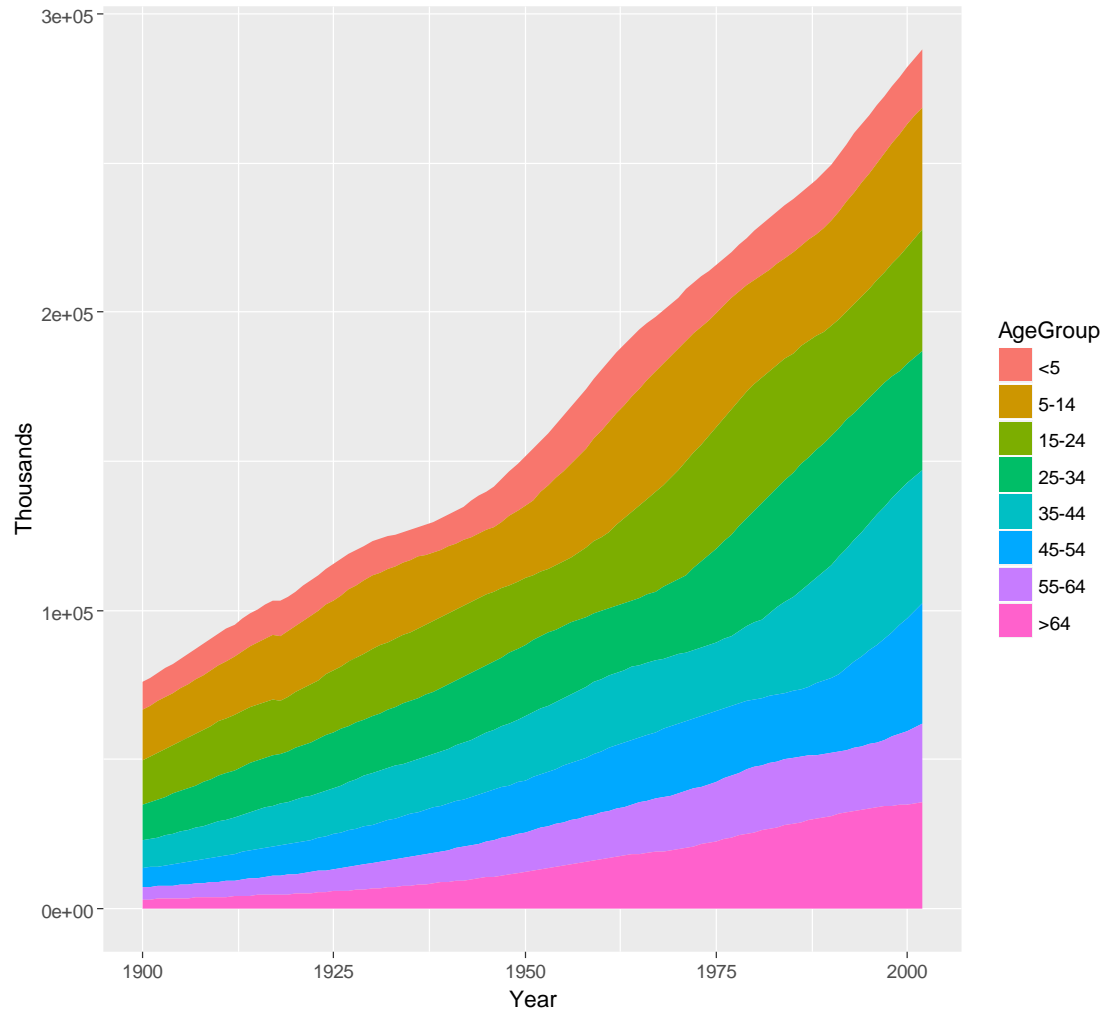
Bag Plot

面积图

- 强调数据随时间而变化的程度，引起对总值趋势的注意



堆积面积图

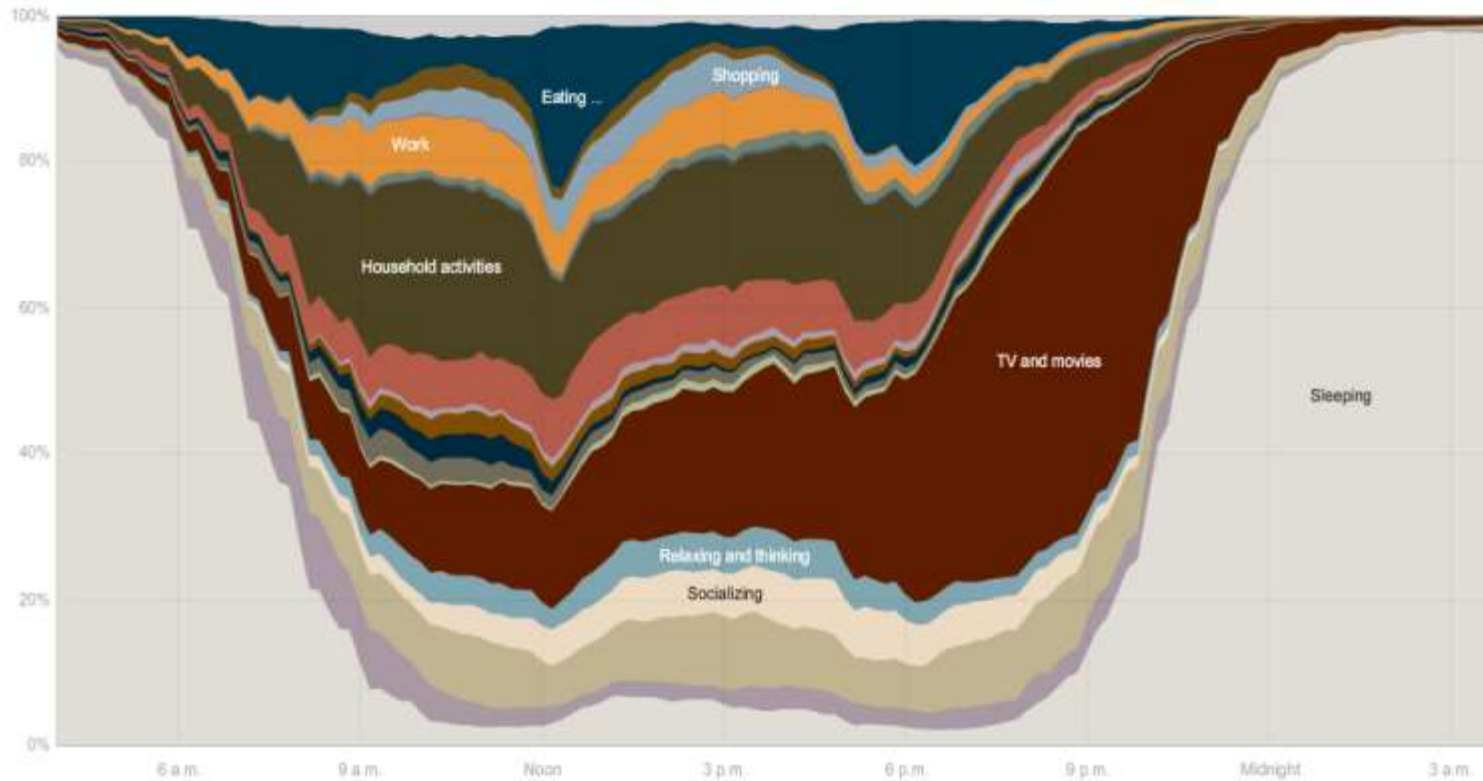


堆积面积图

People ages 65 and over

At 2 p.m., about 1 in 15 people over age 65 is asleep. Older people also spend more time eating (particularly breakfast).

Everyone	Employed	White	Age 15-24	H.S. grads	No children
Men	Unemployed	Black	Age 25-54	Bachelor's	One child
Women	Not in lab.	Hispanic	Age 65+	Advanced	Two+ children

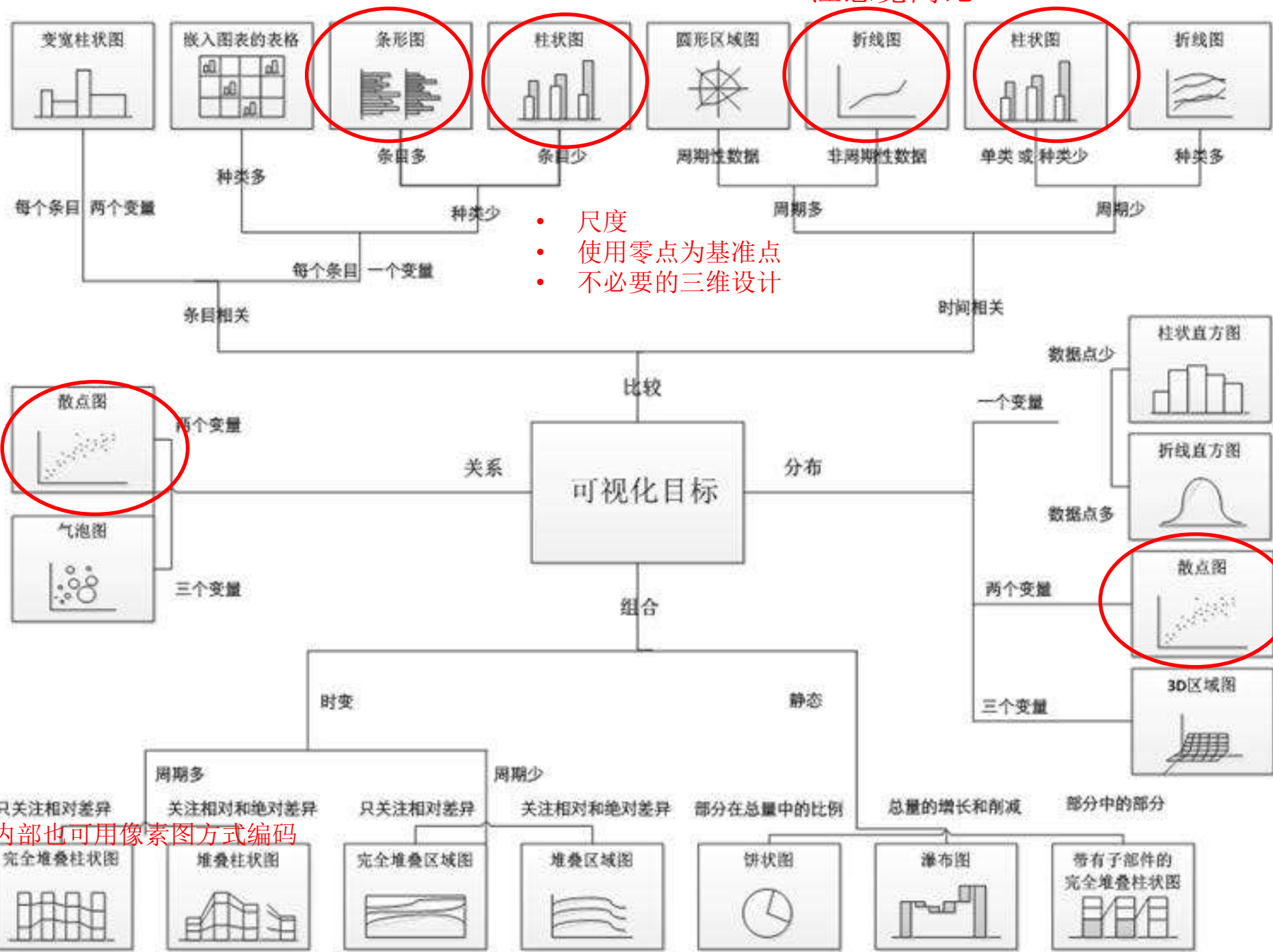


交互的数据可视化

- 交互的可视化可以进一步增强用户对于数据的探索和认知，鼓励用户参与
- D3 (*Interactive Data Visualization for the Web*), Processing.js
- Plotly
- Bokeh...
- <https://scnetworkviz.github.io/SCNetworkViz/>

统计图表

注意宽高比



数据变换的目的是简化数据、降低显示数据的尺寸和复杂度

- 归一化
- 曲线拟合
- 统计采样
- 降维, etc.