

中级应用统计学

导言

张晨峰

华东理工大学商学院

2016年9月13日

参考教材

- 多元统计分析（何晓群，第4版，中国人民大学出版社）
- 概率论及其应用（威廉·费勒，卷1·第3版）
- 统计推断（卡塞拉和贝耶，翻译版·原书第2版）
- 实用多元统计分析（约翰逊和威克恩，第6版）
- 应用多元统计分析（哈得勒和西马，翻译版，第2版）
- 多元数据分析（Joseph F.H, William C.B, etc. 英文版，第7版）

教师信息

联系方式

- 电邮: glen.zhang7@gmail.com
- 手机: 13918610536

课程评分

- 平时成绩: 出勤+平时作业+小项目=30%
- 期末考试: 闭卷考试=70%

课程大纲

- ① 导论
- ② 相关分析和回归分析
- ③ 主成分和因子分析
- ④ 聚类分析
- ⑤ 定性统计分析
- ⑥ 决策分析
- ⑦ 时间序列分析
- ⑧ 判别分析
- ⑨ 结构方程
- ⑩ 大数据和机器学习

课程信息

课程网站

github网站: <https://github.com/plutoese/datascience>

软件程序

- Spss和Stata
- R
- Python

1 导论

主要内容

- 统计学与生活
- 统计学简史
- 应用统计学简介
- 一些应用的实例
- 小项目——词云
- 概率统计简要回顾

1.0 统计学与生活

数据搜集无处不在

You are being WATCHED!

1.1 统计学简史

统计学的历史时刻

公元前450年，伊利斯（Elis）的希皮亚斯利（Hippias）用国王统治时间的均值计算出了第一次奥林匹克运动会的时间，就是在他生活的时代的300年之前。

统计学的历史时刻

中国历史上第一次完整地记载中国各州、郡的户数和人口，是在公元5年（西汉平帝元始五年），据《汉书·地理志》记载，当时有1223.3万户，5959.4万人。

1.1 统计学简史

统计学的历史时刻

公元840年，数学家肯迪（Al-Kindi）应用频谱分析（frequency analysis）进行了密码破解。此外，他还把阿拉伯数字引入到了欧洲。

统计学的历史时刻

公元1570年，天文学家第谷·布拉赫应用算术平均减少了行星观测的误差。

1.1 统计学简史

统计学的历史时刻

公元1657年，物理学家惠更斯（Huygens）发表了第一部概率论著作《论赌博中的计算》（On Reasoning in Games of Chance）。

统计学的历史时刻

公元1805年，勒让德引入最小二乘法对观测数据集进行曲线拟合。

1.1 统计学简史

统计学的历史时刻

公元1894年，卡尔·皮尔逊（Karl Pearson）提出了“标准差”这一术语，并且发展了卡方检验。

统计学的历史时刻

公元1900年，劳伦斯·巴施里耶（Louis Bachelier）将股票价格的涨跌看作是一种随机布朗运动，并据此进行了严格的数学描述，成为了金融数学领域的开创性研究。

1.1 统计学简史

统计学的历史时刻

公元1935年，费雪（R.A. Fisher）出版的《实验设计》（The Design of Experiments）建立了实验设计法的基础，革新了现代统计学。

统计学的历史时刻

公元1979年，布拉德利·埃弗龙（Bradley Efron）提出了自助法（Bootstrap），这是一个很通用的算法，可以计算任意估计的标准误差。

1.1 统计学简史

统计学的历史时刻

公元1997年，大数据（Big Data）这个词首先见诸于出版物中。

统计学的历史时刻

公元2012年，内特·希尔（Nate Silver）成功预测了美国大选中50个州的投票结果，成为了媒体明星。

1.1 统计学简史

统计学的创立时期（17世纪中叶至18世纪中叶）

- **国势学派** 产生于17世纪的德国，其主要代表人物是海尔曼·康令和阿亨华尔。该学派在进行国势比较分析中，偏重事物性质的解释，而不注重数量对比和数量计算，但却为统计学的发展奠定了经济理论基础。
- **政治算术学派** 产生于19世纪中叶的英国，创始人是威廉·配第。他在《政治算术》一书中，利用实际资料，运用数字、重量和尺度等统计方法对英国、法国和荷兰三国的国情国力，作了系统的数量对比分析，从而为统计学的形成和发展奠定了方法论基础。

1.1 统计学简史

统计学的发展时期（18世纪末至19世纪末）

- **数理统计学派** 以比利时的阿道夫·凯特勒为首的统计学家，主张用研究自然科学的方法研究社会现象，正式把古典概率论引进统计学，使统计学进入一个新的发展阶段。
- **社会统计学派** 主要代表人物主要有恩格尔、梅尔等人。他们融合了国势学派与政治算术学派的观点，沿着凯特勒的“基本统计理论”向前发展，但在学科性质上认为统计学是一门社会科学，是研究社会现象变动原因和规律性的实质性科学，以此同数理统计学派通用方法相对立。

1.1 统计学简史

现代统计学的发展时期（20世纪初到现在）

- 统计学的主流从描述统计学转向推断统计学
- 向多分支学科发展
- 统计预测和决策科学的发展
- 信息论、控制论、系统论与统计学的相互渗透和结合
- 计算技术和一系列新技术、新方法在统计领域不断得到开发和应用

1.2 应用统计学简介

统计学的概念

统计学是一门收集数据、分析数据，并根据数据进行推断的艺术和科学。（大英百科全书）

统计学的分支

- 描述统计学
- 推断统计学

1.2 应用统计学简介

描述统计学

致力于数据集的整理、概括以及描述的统计学分支

推断统计学

利用样本数据集对总体作出推断的统计学分支

1.2 应用统计学简介

现代统计分析方法

- **分类分析方法** 聚类分析、判别分析、定性资料分析等
- **结构简化方法** 主成分分析、因子分析、对应分析等
- **相关分析方法** 定性资料分析、典型相关分析、回归分析、主成分分析、因子分析、对应分析等
- **预测决策方法** 回归分析、判别分析、定性资料分析、聚类分析等

1.3 一些应用的实例

红学研究

韦博成教授以情景描写为基础，应用统计学方法研究了《红楼梦》前80回与后40回在若干情景描写上的差异。他选择了花卉、树木、饮食、医药与诗词这5个情景指标，统计出它们在前80回与后40回中出现的频数，并应用统计学中的“等价性检验”方法来检验二者的差异。例如，《红楼梦》在前80回中有34回涉及饮食方面的描写；后40回仅有8回涉及饮食方面的描写。结果表明，《红楼梦》前80回与后40回在饮食与花卉的描写上确实存在非常显著的差异；在树木的描写上也存在明显差异。这提供一个强有力的证据，说明《红楼梦》前80回与后40回在某些文风上确实存在非常显著的差异。

1.3 一些应用的实例

蝎子号事件

1968年5月，美国潜艇蝎子号（Scorpion）在完成北大西洋参观后，在返回纽波特纽斯（Newport News）途中消失了。克雷文组建了一个囊括各方面专家的团队。团队成员包括数学家、潜艇专家和救助人员等。有趣的是，他非但不是要求团队成员互相协商寻求一个答案，反而请每个成员提供自己对每个可能场景的发生概率的猜测。克雷文认为，如果他能把所有答案加在一起，构建一个蝎子号出事全景的复合图像，他应该会得到对潜艇最终位置的很好估计。事实证明这个集体的判断非常精彩。在蝎子号消失后的第五个月，海军发现了它。它和克雷文最后得到的位置只相差约200米。

1.4 小项目——词云



1.4 小项目——词云



1.4 小项目——词云



1.4 小项目——词云



1.4 小项目——词云



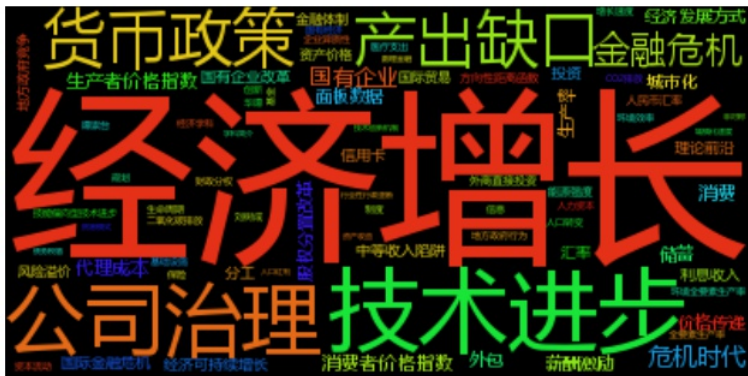
1.4 小项目——词云



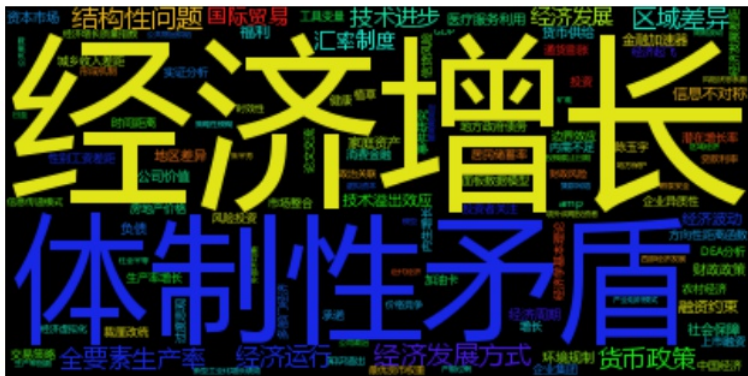
1.4 小项目——词云



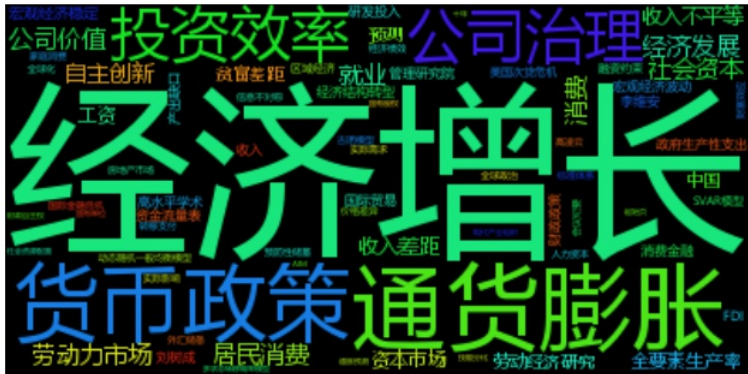
1.4 小项目——词云



1.4 小项目——词云



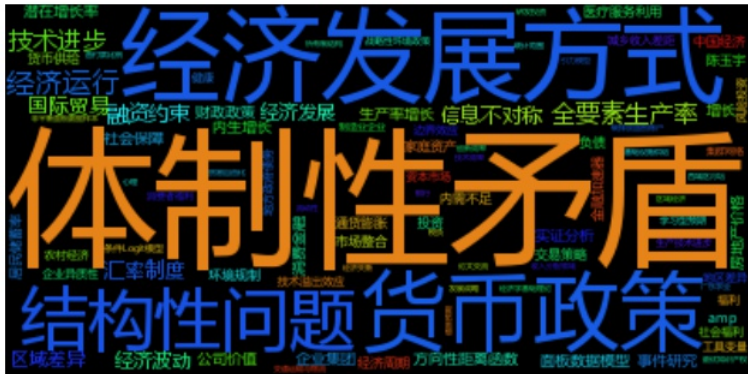
1.4 小项目——词云



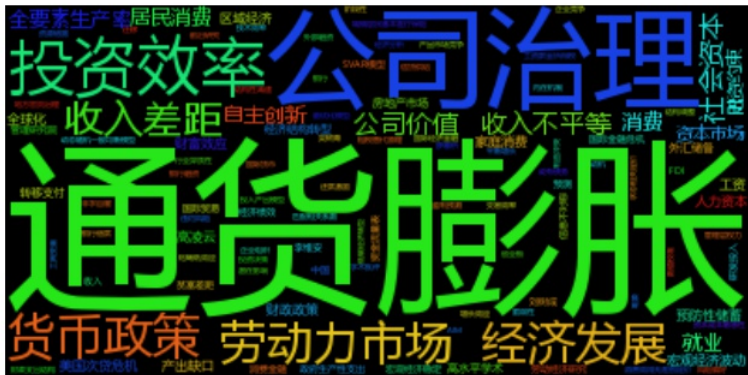
1.4 小项目——词云



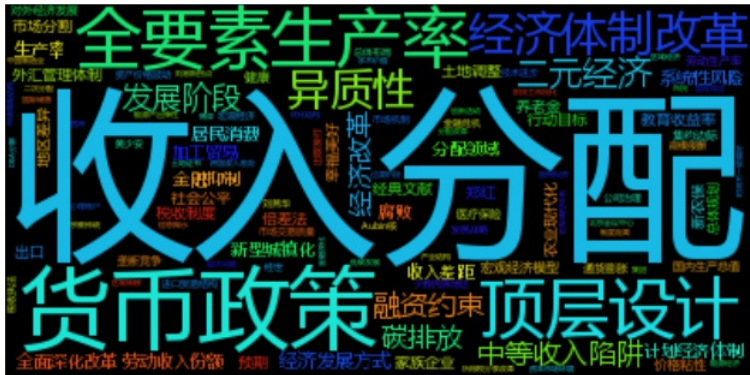
1.4 小项目——词云



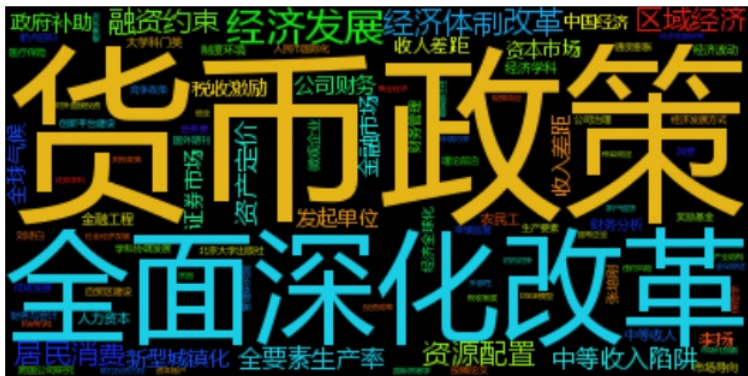
1.4 小项目——词云



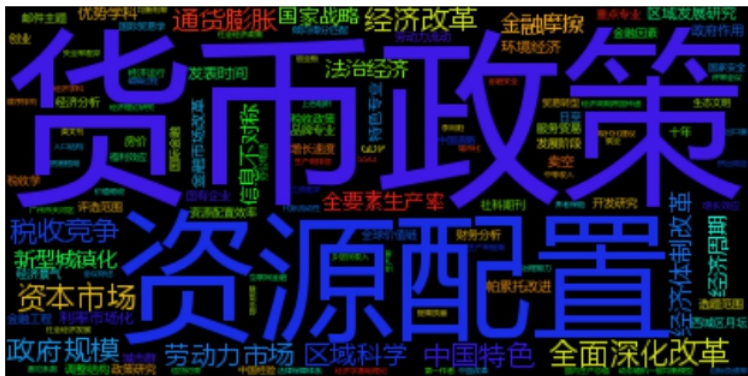
1.4 小项目——词云



1.4 小项目——词云



1.4 小项目——词云



1.5 概率统计简要回顾：集合论

定义 (样本空间)

某次试验全体可能的结果所构成的集合 S 为该试验的样本空间 (*sample space*)。

定义 (事件)

一个事件 (*event*) 是一次试验若干可能的结果所构成的集合，即 S 的一个子集 (可以是 S 本身)

集合的关系和运算

集合的关系 序关系，相等关系

集合的运算 并、交、补

1.5 概率统计简要回顾：集合论

定义 (不交)

如果两个事件 $A \cap B = \emptyset$ ，称两个事件 A 和 B 不交 (*disjoint*)

定义 (两两不交)

如果对于任意 $i \neq j$ ，都有 $A_i \cap A_j = \emptyset$ ，则称事件 A_1, A_2, \dots 两两不交 (*pairwise disjoint*)。

定义 (划分)

如果事件 A_1, A_2, \dots 两两不交，并且 $\bigcup_{i=1}^{\infty} A_i = S$ ，则称 A_1, A_2, \dots 构成 S 的一个划分 (*partition*)。

1.5 概率统计简要回顾：概率论的公理化基础

定义 (σ 代数)

S 的一族子集如果满足下列三个性质，就称为一个 σ 代数，记作 Λ 。

- $\emptyset \in \Lambda$
- 若 $A \in \Lambda$, 则 $A^c \in \Lambda$
- 若 $A_1, A_2, \dots \in \Lambda$, 则 $\bigcup_{i=1}^{\infty} A_i \in \Lambda$

1.5 概率统计简要回顾：概率论的公理化基础

定义 (概率函数)

已知样本空间 S 和 σ 代数 Λ ，定义在 Λ 上且满足下列条件的函数 P 称为一个概率函数 (*probability function*)

- 对任意 $A \in \Lambda$, $P(A) \geq 0$
- $P(S) = 1$
- 若 $A_1, A_2, \dots \in \Lambda$ 且两两不交, 则 $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

1.5 概率统计简要回顾：条件概率和独立性

定义 (条件概率)

设 A, B 为 S 中的事件, 且 $P(B) > 0$, 则在事件 B 发生的条件下事件 A 发生的条件概率记作 $P(A | B)$, 表示为

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

定义 (统计独立)

称事件 A, B 统计独立, 如果

$$P(A \cap B) = P(A)P(B)$$

1.5 概率统计简要回顾：随机变量

定义 (随机变量)

从样本空间映射到实数的函数称为随机变量 (*random variable*)

定义 (累积分布函数)

随机变量 X 的累积分布函数 (*cumulative distribution function*, 简记为*cdf*), 记作 $F_X(x)$, 表示

$$F_X(x) = P_X(X \leq x)$$

其中 x 任意。

1.5 概率统计简要回顾：随机变量

定理 (累积分布函数的性质)

函数 $F_X(x)$ 是一个累积分布函数，当且仅当它同时满足下列三个条件：

- $\lim_{x \rightarrow -\infty} F(x) = 0$ ，且 $\lim_{x \rightarrow +\infty} F(x) = 1$
- $F(x)$ 是 x 的单调递增函数
- $F(x)$ 右连续

定义 (连续和离散)

设 X 为一随机变量，如果 $F_X(x)$ 是 x 的连续函数，则称 X 是连续的 (*continuous*)；如果 $F_X(x)$ 是 x 的阶梯函数，则称 X 是离散的 (*discrete*)。

1.5 概率统计简要回顾：随机变量

定义 (概率质量函数)

离散随机变量 X 的概率质量函数 (*probability mass function*, 简称*pmf*) 为

$$f_X(x) = P_X(X = x)$$

定义 (概率密度函数)

连续随机变量 X 的概率密度函数 $f_X(x)$ (*probability density function*, 简称*pdf*) 是满足下式的函数

$$F_X(x) = \int_{-\infty}^x f_X(t) dt$$

1.5 概率统计简要回顾：期望

定义 (期望)

连续随机变量 $g(X)$ 的期望 (*expected value*) 或均值 (*mean*), 记作 $Eg(x)$, 定义为

$$Eg(X) = \int_{-\infty}^{+\infty} g(x)f_X(x) dx$$

1.5 概率统计简要回顾：期望

定义 (矩)

对任意整数 n ， X （或 $F_X(x)$ ）的 n 阶矩，记作 μ'_n ，定义为

$$\mu'_n = EX^n$$

X 的 n 阶中心矩，记为 μ_n ，定义为

$$\mu_n = E(X - \mu)^n$$

1.5 概率统计简要回顾：多维随机变量

定义 (联合概率密度函数)

设 (X, Y) 为连续随机变量，称 R^2 到 R 的函数 $f(x, y)$ 为 (X, Y) 的联合概率密度函数 (*joint pdf*)，如果对于任意 $A \in R$ ，都有

$$P((X, Y) \in A) = \iint_A f(x, y) dx dy$$

X 和 Y 的边缘概率密度函数分别为

$$f_X(x) = \int_{-\infty}^{+\infty} f(x, y) dy$$

$$f_Y(y) = \int_{-\infty}^{+\infty} f(x, y) dx$$

1.5 概率统计简要回顾：多维随机变量

定义 (独立随机变量)

设 (X, Y) 是二维随机变量，其联合概率密度函数为 $f(x, y)$ ，边缘概率密度函数分别为 $f_X(x)$ 和 $f_Y(y)$ 。称 X 和 Y 是独立随机变量，如果对于任意 $x, y \in R$ ，都有

$$f(x, y) = f_X(x)f_Y(y)$$

定义 (协方差)

随机变量 X 和 Y 的协方差定义为：

$$\text{Cov}(X, Y) = E((X - \mu_X)(Y - \mu_Y))$$

1.5 概率统计简要回顾：随机样本的性质

定义 (随机样本)

如果随机变量 X_1, \dots, X_n 相互独立且有相同的边缘概率密度函数 $f(x)$ ，则称 X_1, \dots, X_n 是总体 $f(x)$ 的大小为 n 的随机样本，或称 X_1, \dots, X_n 是概率密度函数为 $f(x)$ 的独立同分布随机变量（简记为*iid*随机变量）。

1.5 概率统计简要回顾：随机样本的性质

定义 (统计量)

设 X_1, \dots, X_n 是从总体中抽取的大小为 n 的随机样本， $T(x_1, \dots, x_n)$ 是定义在 (X_1, \dots, X_n) 的样本空间上的实值或向量值函数，则随机变量或随机向量 $Y = T(X_1, \dots, X_n)$ 称为一个统计量 (*statistic*)， Y 的概率分布称为 Y 的抽样分布 (*sampling distribution*)。

1.5 概率统计简要回顾：随机样本的性质

定义 (样本均值)

样本均值 (*sample mean*) 是随机样本值的算术平均，常记作

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i$$