

# 统计

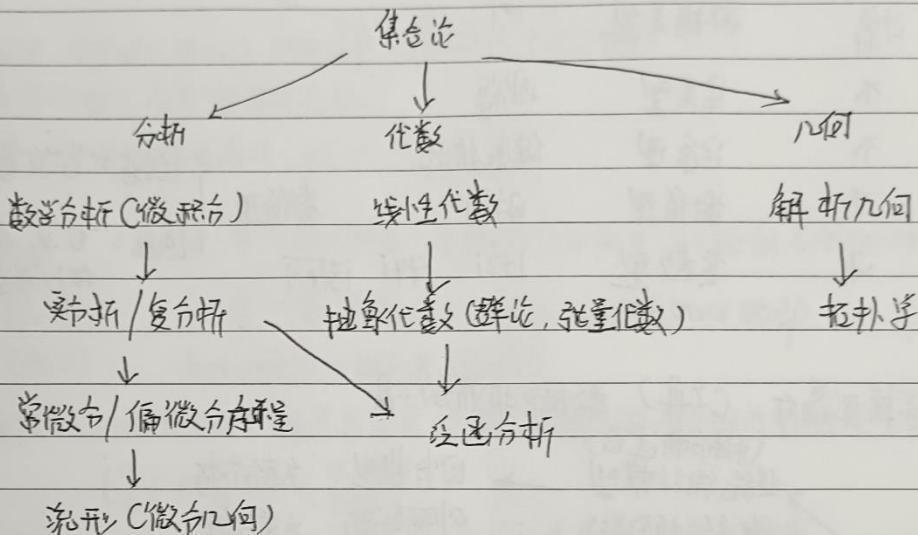
第一天：预备知识，线性函数

第二天：函数、微积分

第三天：数据度量，抽样分布，参数估计，假设检验

第四天：列联表分析，相关分析，回归分析

## 第一节：数学概况



## 第二节：数据类型

### 离散型数据

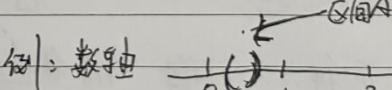
离散随机变量是指一个只取有限个或可数无限个数值的随机变量。

通常用古典概型来描述。

### 连续型数据

连续随机变量是指一个取任何实数的概率皆为零的变量。  $\frac{1}{\infty} = 0$  是没办法通过数

通常用几何概型来描述。

图：数轴  在数轴上取一个区间   
 整数  $N$ (全体自然数) 如果区间内有可能不包含这个数集里的任意一个元素, 那么此数集为离散数集

连续  $Q$ (全体有理数) 无论区间取在哪, 无论区间大小, 此区间都至少含有这个数集里的元素, 那么这个数集为稠密数集

连续  $R$ (全体实数) 无论区间取在哪, 无论区间大小, 此区间都至少含有这个数集里的元素, 且此数集是连续的, 那么这个数集是连续的

离散型数据直接数频数就可以计算出来,

数据一但带上时间属性就天然带上了某种相关性(过去与未来的相关性)

**数据类型**

**时间序列数据**: 在一个时间段或一个时间点内取到的所有数据叫做时间序列数据  
时间序列数据: 随着时间会变化的数据  
**面板数据**: 具有时间属性, 又有空间属性的数据, 即时间属性  
例: 去年企业在各个省市的销量  
空间属性

排序	计算	数据类型	例
不	不	定量型	固定
同	不	定性型	健康状况
同	可	数值型	时间
不	可	复数型	1+2i 2+i i=1

**第二步** 数据模型设计 (工具) 数据分析的工具  
 (甲壳虫模型去)  
 传统统计模型 → 回归模型, 线性分析,  
 (解决结构问题) 列联表分析, 方差分析,  
 等等等等...  
**数据挖掘模型**  
 (消费者行为)  
 数据挖掘模型 → 大数据, 神经网络,  
 (解决预测问题) 支持向量机, 决策树网络,  
 关联规则, 等等等等...

**支持向量机**  
 例:

支持向量机: 基本上为二分类模型  
 多个分类模型拼接在一起可以解决多类问题  
 例: 人脸识别: 把中国人根据脸部分割成4亿类。

**现代数据分析**  
 但是这些往往差别很大。  
 把数据放入一个数据仓库中, 通过对数据仓库的各种形状与分析从而挖掘出数据  
 最底层所携带的信息。

**线性化函数**

**1. 向量**

物理	统计学	计算机
有方向的量	空间中的点	坐标

在空间里, 我们还可以将以上两种向量的相加合到一起。  
 当我们默认向量的起始在原点  
 向量一般用小写字母表示, 如  $a, \beta, u, v$ 。

**向量的坐标**: 表示中, 第一个数(向量)告诉我们在原点沿x轴的方向如何移动。  
 $[2, 3] [1, 2]$  第二个数告诉我们从原点沿y轴如何移动。  
 $[2, -1] [4, 0]$  此为二维空间的向量的坐标  
 向量的加法: 2维空间内, 就是连接2个向量所围成的平行四边形的对角线。  
 $(a_i + b_i) = (a_i + b_i)$   
 向量的减法: 将向量向量按四边形法则(拉伸), 反数表示反向操作。  
 $a(-i) = (-a)$   
**补充: 向量的乘法**

**线性组合与向量空间**

**线性组合**: 将一个向量组中的向量做线性运算相加, 即得到该向量组的一个所谓的线性组合。  
**定义如下**: 空间 V 中的一组向量  $v_1, \dots, v_n$  的线性组合是形如  
 $\alpha_1 v_1 + \dots + \alpha_n v_n$  的向量, 其中  $\alpha_1, \dots, \alpha_n$  是常数。  
**例**: 在 R<sup>n</sup> 空间中,  $(1, 1) + (2, 3) + (2, -1)$  的线性组合, 因为  $(\frac{1}{3}) + 2(\frac{2}{3}) + 2(-\frac{1}{3}) = (\frac{1}{3})$   
**生成空间**: (以二维空间为例) 在线性组合中, 给定两个非零向量, 一个向量固定, 另一个向量自由变化, 其线性组合将得到一条直线。  
 但是, 两个向量都固定时, 可以得到一个平面吗?  
 大多数情况下确定如此, 除非两向量共线。

**统计学分析数据的方法**

**描述性分析**: 研究数据收集、处理和描述的  
 总体指标: 对比差异、集中趋势、离散程度、偏态、峰态, ...

**推断性分析**: 研究如何利用样本数据来推断总体特征的统计方法  
 估计, 假设检验, 判决分析, 方差分析, 相关分析, 回归分析, ...

**统计学的基本概念**、**一元模型**

**结构化数据**: 指所有可以用电子表示的  
 非结构化数据 转成 结构化  
 再进行分析

1. 与山羊的数据是数据  
 数字(结构化数据) 可以进行比较, 加法乘法等运算, 严格的数据格式  
 2. 数据的形式  
 文本(非结构化数据) 不可运算, 如名字等  
 正则表达式(万能表达式) → 词模型

只要能辨别信息的数据才是数据。注: 没有信息也是信息  
 东西 例如钱的东西对于我们来说是无关的

**数据的分类**

分类型数据: 对事物进行分类的, 如人的性别男/女, 喜欢/不喜欢, ...  
 1. 按照计量尺度分类  
 顺序型数据: 对事物类别排序的, 如年龄: 一岁, 三岁, 三岁半  
 整体型数据: 对事物的精确度, 如身高: 175cm, 180cm

2. 特点:  
 分类型数据: 不可排序, 不可计算  
 顺序型数据: 可排序, 不可计算  
 整体型数据: 可排序, 可计算  
 PS:  
 整体型数据: 不可排序, 不可计算  
 例: 1+2i; 2+3i; i=1  
 顺序型数据: (1) 整数首尾相加  
 查数单位: 整数的极限收敛于一个单位

**数据挖掘**, paperweekly, arxiv.org (前沿文章), 中国国家数据中心

**数据的其他分类**

1. 根本原因  
 原始数据 (一手数据, 原始资料)  
 衍生数据 (二手数据, 次级资料)

2. 按数据的不同  
 实验的数据  
 (微观数据): 在一个时间点或一个时间段内分析数据

3. 按与时间的关系  
 与时间的数据  
 (宏观数据):  
 1. 元数据(元数据) 产生操作的内部给程序员的数据  
 2. 生产数据(生产性的数据) 生成模型的数据  
 3. 交易数据(日常流水) 指向用户  
 按时间  
 企业 → 端 → 系统 → 表 → 服务 → 标准 → 数据仓库 → 算法  
 1. 物流系统  
 2. 生产系统  
 3. 生产系统  
 主数据一般从生产系统中取  
 数据模型通过数据仓库的连接

**构建数据仓库流程**

ps. 时间序列数据分析与争议

构建  
 数据仓库  
 1. ETL  
 2. SQL  
 3. 自动化 python (层出不穷)

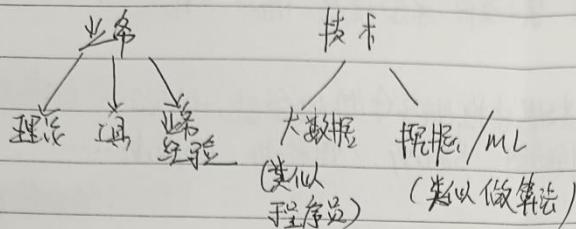
时间序列  
 统计方法  
 人工智能分析方法 (LSTM)  
 例: ARIMA (传统) → 宏观经济  
 商业问题中 常规用 ARIMA  
 但如果用 LSTM 会自动调整机器学习模型  
 LSTM (长短时记忆网络)

**统计学的基本概念**、**一元模型**

# SQL 熟练度

统计的基本概念

推进、推断、回归、逻辑



统计学定义

统计学是一门收集、处理、分析、解释数据并从数据中得出结论的科学

数据采集  
爬虫  
数据接入

数据采集初步步骤  
核心：数据

收集数据 → 处理数据 → 分析数据 → 解释数据

↓  
描述性统计分析      推断性统计分析

数据埋点：使用在网站上的行为数据 (根据业务解进行埋点)

收集  
数据  
阶段

取出流水线上关键点的数据  
爬虫：(爬虫具有法律风险)

数据接入：从各个不同的系统中去取数，同时可以接入外部系统 (数据统计口径和标准)  
最后接出的应该是 - 张表，如果是三张表 就没有什么意思

例：填充缺失值

如果是对行数据的处理，一般称为 数据清洗

如果是对列数据的处理，一般称为 特征工程

处理数据

一列为一字段    一列为一变量    一列为一特征



一列为一“维度”(维度)

分布数据

特征工程用的树可以优化算法



解释数据：某些数据是有趣的语法