

# 数据分析与可视化技术

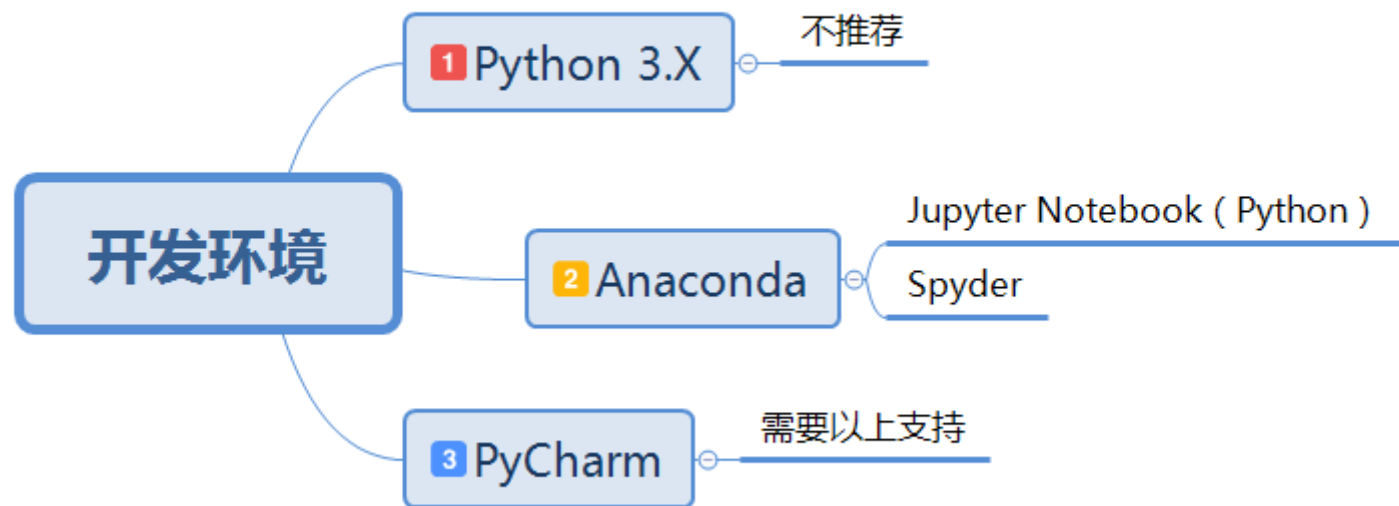
## Data Analysis for Infomatics

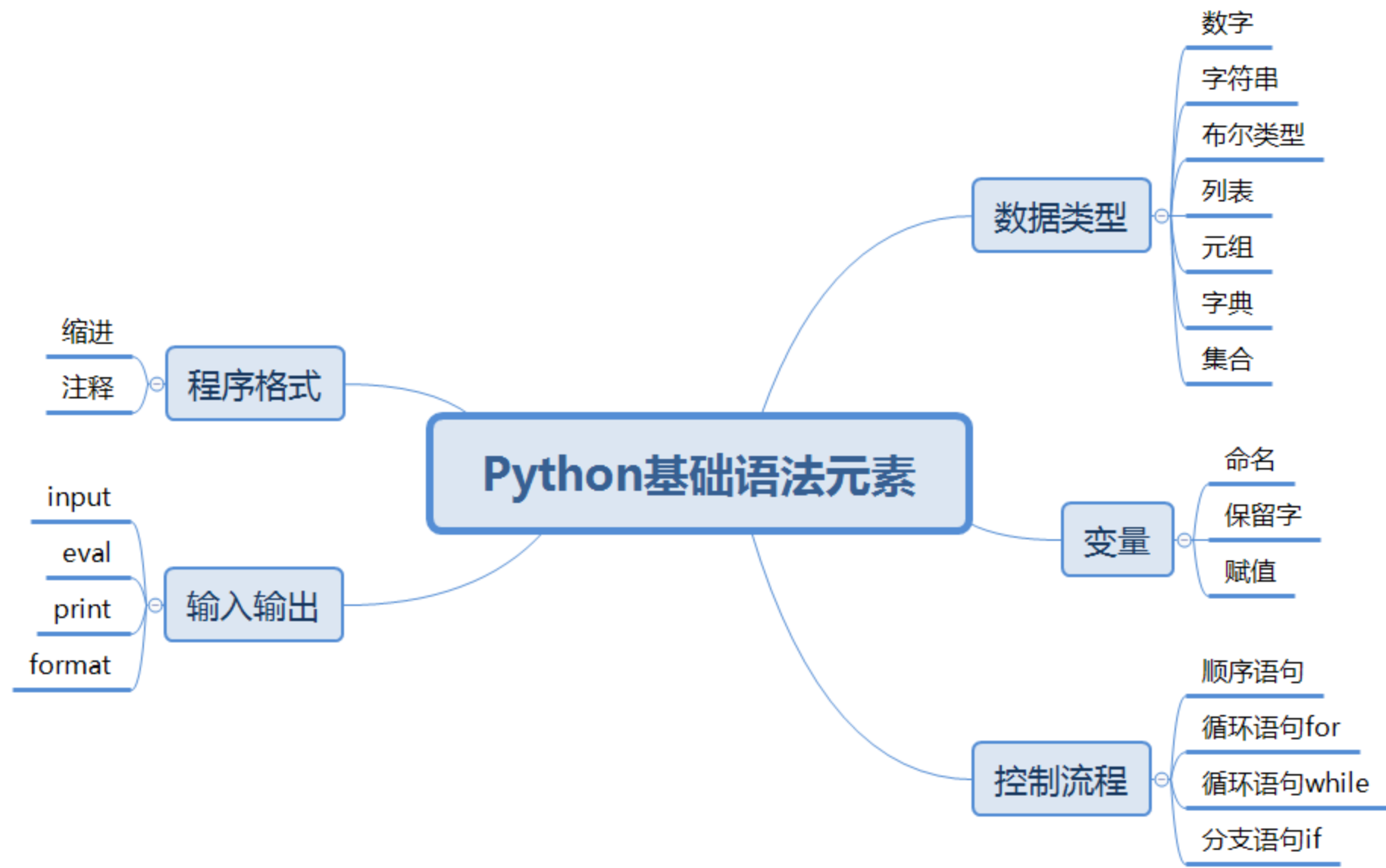
### 第一讲 Python提升探索, Numpy & Pandas

刘子瑜 副研究馆员

[liuziyu@cau.edu.cn](mailto:liuziyu@cau.edu.cn) 中国农业大学图书馆情报研究中心

[https://github.com/CAU-LiuZiyu/Data\\_Analysis\\_for\\_Infomatics](https://github.com/CAU-LiuZiyu/Data_Analysis_for_Infomatics)





# Python 基础数据类型

## 数值类型

分类 • int float complex

操作符 • + - \* / % // \*\*

操作函数 • abs() max() min()

## 字符串类型

基本性质 • 索引 切片

操作符 • + \* in

操作函数 • len() count() strip()

## 布尔类型

逻辑运算 • and or not

## 类型转换

类型判断 • type() isinstance()

类型转换 • str() int() float()

# Python 组合数据类型

## 列表

生成列表 • 直接创建 list()

列表性质 • 索引切片

列表元素的增、删、查、改

操作列表 •

列表的复制

列表的排序、翻转

遍历列表

## 元组

元组性质 • 不可变的列表

操作元组 • 不支持任何关于元组本身的修改

常见用法 • 打包和解包

## 字典

生成字典 • 键的要求

字典性质 • 索引

操作字典 •

键值的添加、删除、修改

字典的遍历

## 集合

集合表达 • 没有值的字典

集合运算 • 交集、并集等

操作集合 • 元素的增加、删除、遍历

# Python 程序控制结构

## 条件测试

True or False

比较运算 • > < >= <= == !=

逻辑运算 • and or not

存在运算 • s in S

## 分支结构if

单分支 • if

二分支 • if else

多分支 • if elif else

嵌套语句 • if if if

## 遍历循环for

迭代对象 • 列表、元组、集合、字符串

循环控制 • range()

for else

## 注意事项

尽可能减少嵌套

避免死循环

封装过于复杂的判断条件

## 无限循环while

循环控制

break

continue

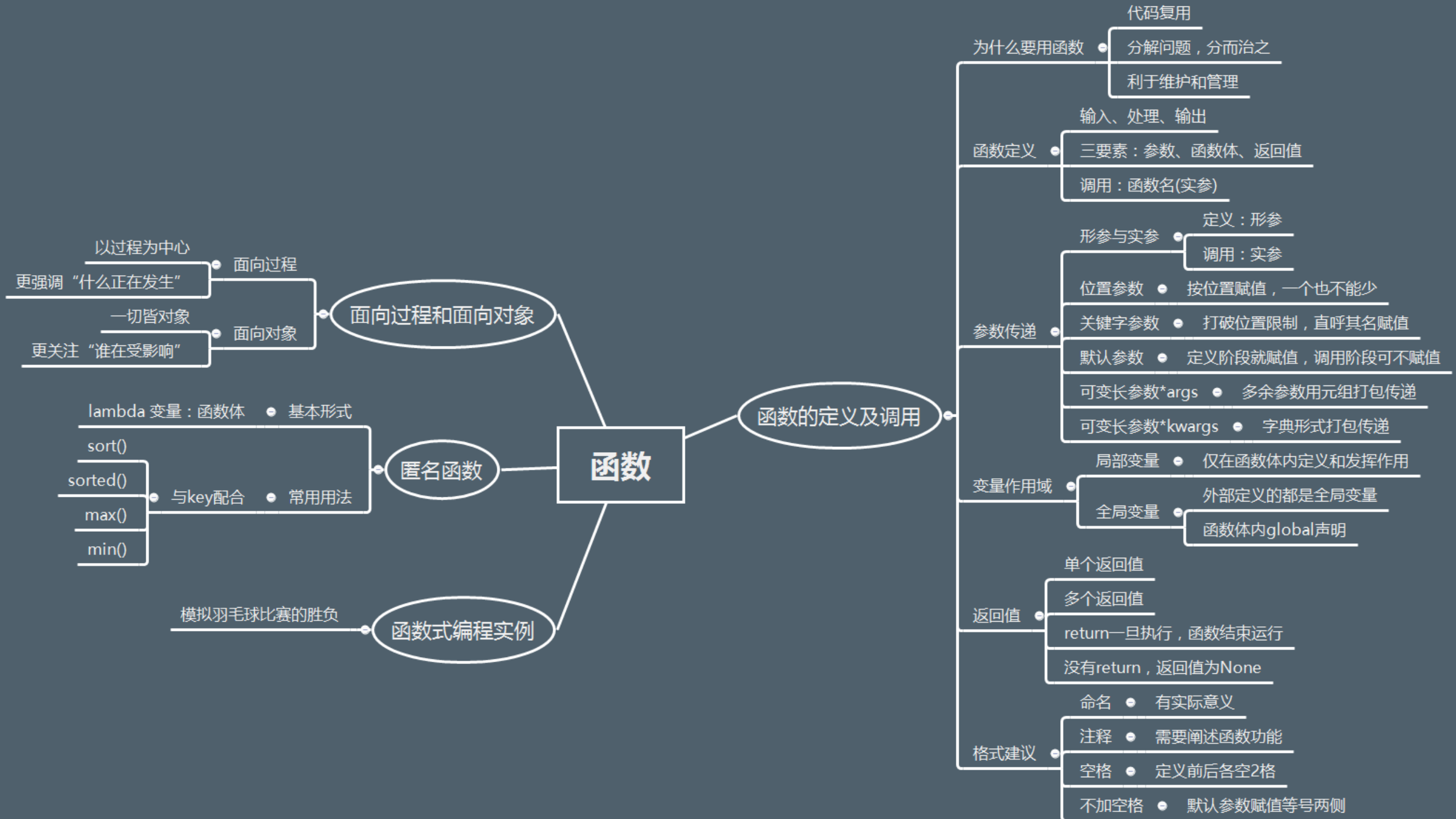
how

风向标

while else

更多例子

没有值的字典 • why



# 类与对象

## 为何要创建对象

why

符合对客观世界的抽象和理解

一切皆对象

一切对象都有自己的内在属性

一切行为都是对象的行为

how

类是对象的载体

把一类对象的公共特征抽象出来，创建通用的类

## 类的定义

类的命名

有实际意义

驼峰体

ElectricCar

组成的单词首字母大写

类的属性

类内部的变量

def \_\_init\_\_ 初始化

self.v = v

类的方法

类内部的函数

## 创建实例

创建实例

实例名+类名（必要的初始化参数）

my\_new\_car = Car('BYD', 'Song Pro', 2020)

访问属性

实例名.属性名

my\_new\_car.brand

调用方法

实例名.方法名（必要的参数）

my\_new\_car.get\_info()

修改属性

通过属性的调用直接修改

通过方法的调用间接修改

## 类的继承

继承的本质

底层抽象继承高层抽象

随着继承，公共特性逐渐增加

简单的继承

自动继承所有方法

先继承

后定义

给子类添加属性和方法

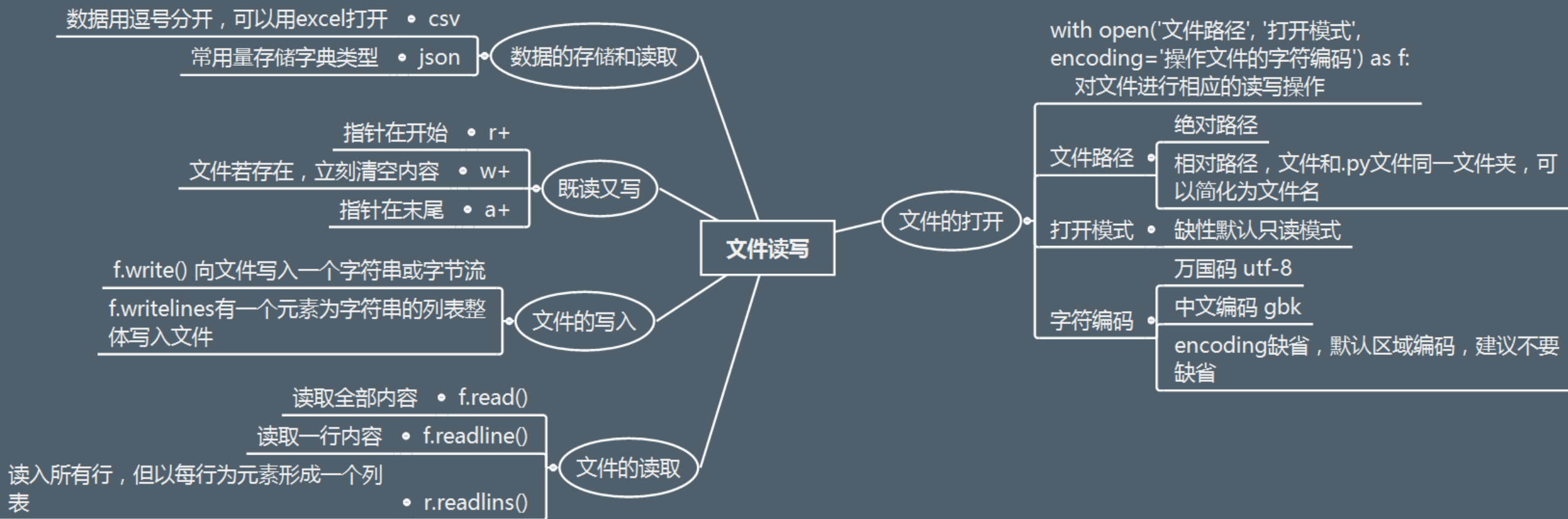
重新父类方法

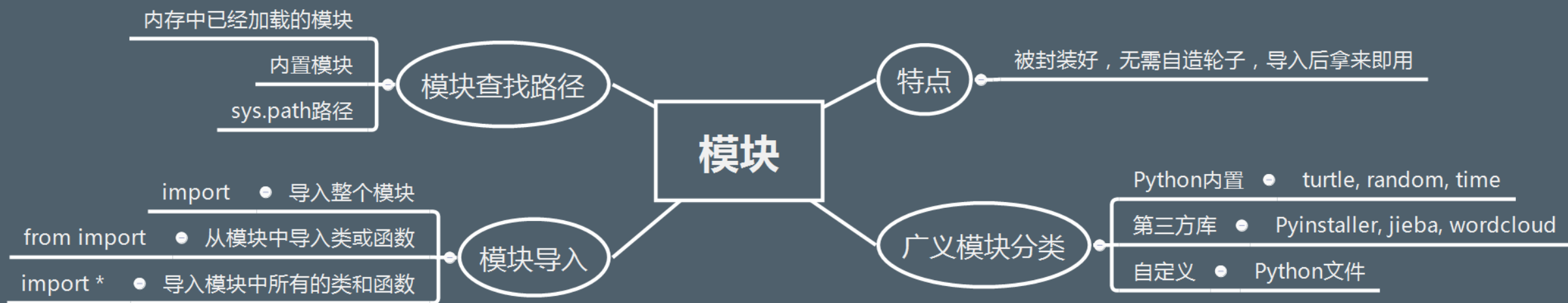
方法名称不变，重写内容

优先调用子类的方法

super().\_\_init\_\_() • 声明继承父类的属性







# Python标准库

## itertools

- 排列组合
  - product 笛卡尔积
  - permutations 排列
  - combinations 组合
  - combinations\_with\_replacement 元素可重复组合
- 拉链
  - zip 短拉链
  - zip\_longest 长拉链
- 无穷
  - count 计数
  - cycle 循环
  - repeat 重复
- 其他
  - chain 锁链
  - enumerate 枚举
  - groupby 分组

## collections

- 具名元组
  - namedtuple 元组的子类
- 计数器
  - Counter 字典的子类
  - most\_common
  - element
- 双端队列
  - deque
  - 两侧进行高效的插入和删除操作

## time

- 获取当前时间
  - localtime 本地时间
  - gmtime 世界统一时间
- 时间戳和计时器
  - time() 纪元为计时起点
  - perf\_counter() 任取一点为计时起点, 记录sleep
  - process\_time() 任取一点为计时起点, 不记录sleep
  - 两个时间戳之差可以统计时间差
- 格式化输出
  - strftime()
- 睡觉
  - sleep()

## random

- 伪随机数, 可满足大多数工程需要
- 随机种子
  - seed
    - 相同种子产生相同随机数
    - 不设定种子, 默认为系统当前时间
- 随机整数
  - randint(a, b) [a, b)
  - randrange(a) [0, a)
  - randrange(a, b, step) [a, b), 以step为步长
- 随机浮点数
  - random() [0.0, 1.0)
  - uniform(a, b) [a, b)
- 用于序列的函数
  - choice() 返回1个
  - choices(seq, weights=None, k) 重复采样
  - shuffle(seq) 重新排列
  - sample(pop, k) 不重复采样
- 概率分布
  - 产生符合相应分布的数

# 程序异常处理

## 异常的处理

单分支

多分支

• try except

万能异常Exception

捕获异常的值 as

try顺利执行，则else也执行 • try\_except\_else

无论try执行与否，finally都执行 • try\_except\_finally

## 三类Bug

语法错误

语义错误

除零运算——ZeroDivisionError

找不到可读文件——FileNotFoundError

值错误——ValueError

异常 • 类型错误——TypeError

NameError

KeyError

...

当异常发生时，如果不预先设定处理方法，程序就会中断

# Python提升探索

## 三大神器

### 生成器

无需海量存储，边执行边计算

生成器表达式

生成器函数yield

### 可迭代对象

列表、元组、字符串等

生成器

列表等不是迭代器，可通过iter()创建

生成器都是迭代器

itertools工具

文件时迭代器

range()不是迭代器，是懒序列

迭代器被next()函数调用，直至耗尽

### 迭代器

### 迭代器

### 装饰器

修改函数功能，但不修改源代码和接口

函数对象、高阶函数

嵌套函数、闭包

简单的装饰器、装饰带参函数、带参数的装饰器

立即执行装饰器、用wraps回复被装饰函数本源

## 数据类型的底层实现

### 列表

元素分散存储在内存中

存储的是元素地址

元素地址是连续存储的

存储与一个稀疏数组

### 字典

通过键的散列值确定存储位置

### 字符串

通过紧凑数组直接存放数据

### 是否可变

可变：id不变，内容不变

不可变：内容一变，id就变了

## 简洁的语法

### 解析语法

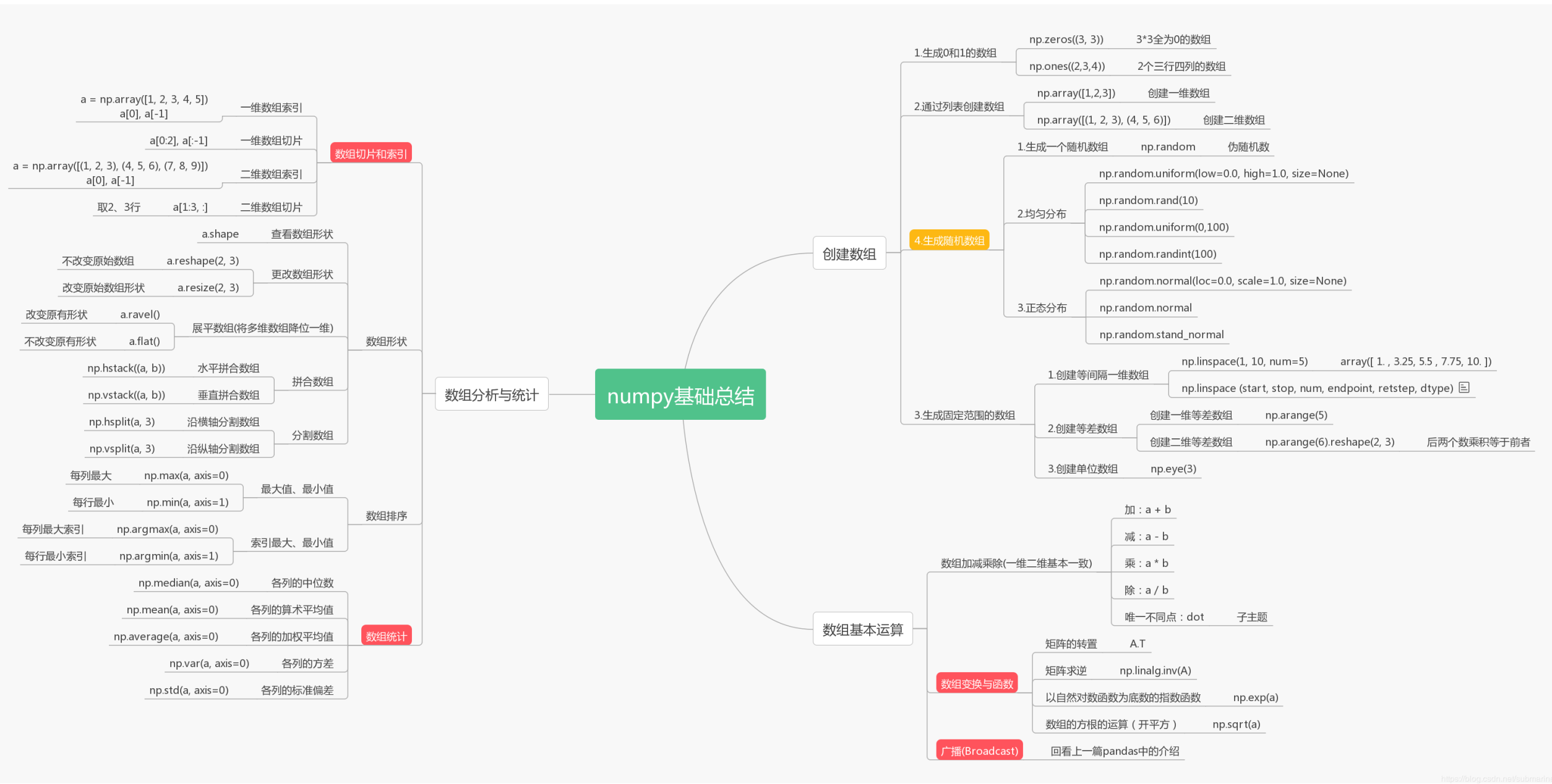
[expression for value in iterable if condition]

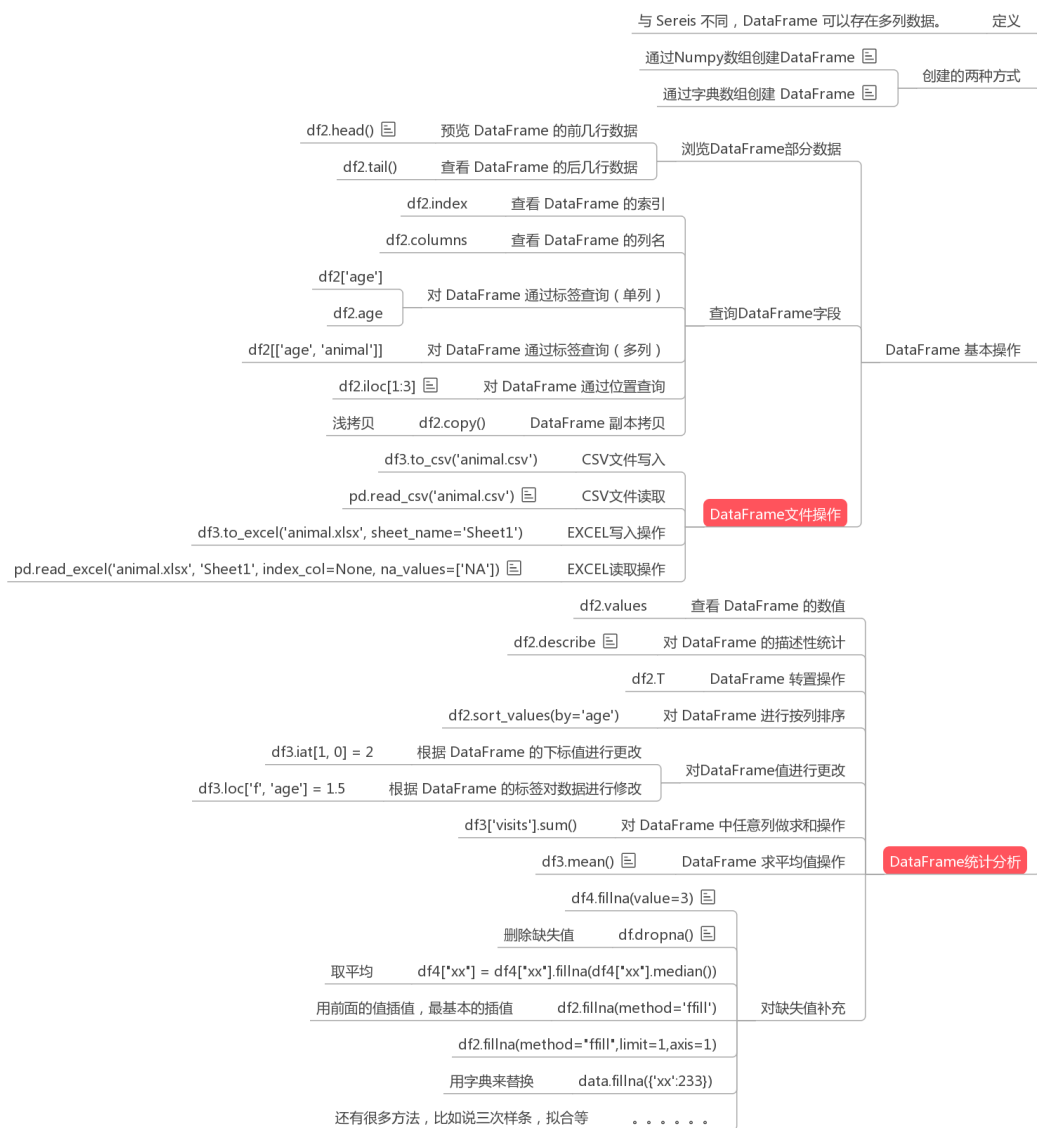
表达式

支持列表、字典、集合、生成器推导

### 条件表达式

expr1 if condition else expr2





DataFrame

## pandas总结

Series

