# 数据挖掘和大数据分析



## **Outline**



1 Review Decision Tree (Assignments)

2 Association Rules



3 Apriori Algorithm



### 作业清单 (5/20) -

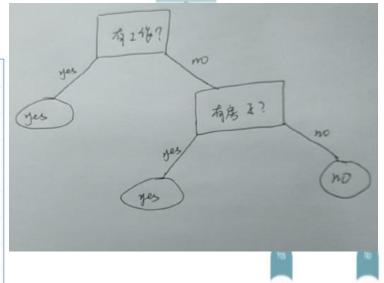
【1】 数据<u>集如下</u>图所示,根据我们对决策树的理解,设计一棵决策树,并输入{Age:36,Salary:H,STU:No,Credit:OK}测试数据,是否与预期结果一致? @注意,不允许直接调用Sklearn提供的决策树方法。

Age #	Salary⊭	STU⊭	Credit <b>₽</b>	Buy Computer
<30₽	He	No₽	OK₽	No₽
<30 ₽	Н₽	No₽	Good₽	No₽
30-40 ₽	Н₽	No₽	OK₽	Yes₽

【2】数据<u>集如下</u>图所示,计算这棵决策树的类别信息熵,并计算基于每个特征值的类别信息熵。根据公式计算每个特征值的信息增益, 画出一棵决策树 (用签字笔画出即可)。。

	ID	年龄	有工作	有自己的房子	信贷情况	类别(是否个给贷款)
	1	青年	否	否	一般	否
	2	青年	否	否	好	否
	3	青年	是	否	好	是
	4	青年	是	是	一般	是
	5	青年	否	否	一般	否
	6	中年	否	否	一般	否
	7	中年	否	否	好	否
	8	中年	是	是	好	是
	9	中年	否	是	非常好	是
	10	中年	否	是	非常好	是
	11	老年	否	是	非常好	是





知识点: 信息熵、决策树、信息增益、条件信息熵

【3】 根据下述数据集,编写代码实现信息熵、条件熵、信息增益等

决策树的关键环节,并撰写实验报告。(最后一列是类别:是否提供贷款)↓

dataSet = [0, 0, 0, 0, 'no'],

#数据集↓

[0, 0, 0, 1, 'no'],

[0, 1, 0, 1, 'yes'],

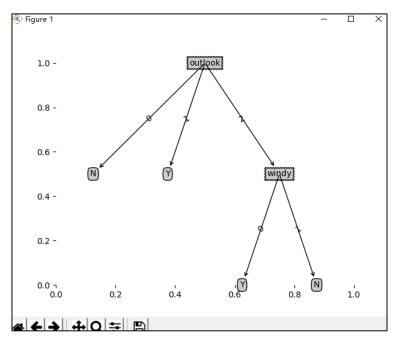
[0, 1, 1, 0, 'yes'],

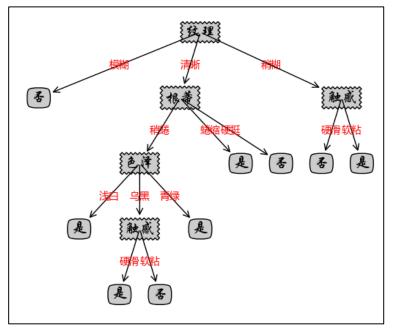
[0, 0, 0, 0, 'no'],

[1, 0, 0, 0, 'no'],

[1, 0, 0, 1, 'no'],

【4】对问题【3】中实现的决策树,实现可视化。(选做)



















Python语言是一种面向对象、动态数据类型的解释型语言,是数据分析师首选编程语言之一.

#### 符号标记

 x | 一个变量,其值为2.
 s | 字符串 (string) 对象.

 l,L | 列表 (list) 对象.
 t | 元组 (tuple) 对象.

 e,E | 集合 (set) 对象.
 d | 字典 (dict) 对象.

#### 基本操作

x = 2 | 定义一个新的变量x, 其值为2. print | 打印输出. # | 行內注释. \*\*\* | 行內注释. \*\*\*\* | 内省,显示对象的通用信息. \*\*\* | 內省,显示对象的通用信息. \*\*\*\* | 內省,显示一对象的那功文档. \*\*\* | 小店们,以下一个对象的帮助文档. \*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*\*\* | \*\*

#### 数据类型及相互转换

$$\label{eq:type} \begin{split} & \text{type}(\mathbf{x}) \mid \pm \mathbf{a} \mp \mathbf{y} \\ & \text{int}(\mathbf{x}) \mid 将 \mathbf{y} \\ & \text{int}(\mathbf{x}) \mid 将 \mathbf{y} \\ & \text{son} \\ & \text{str}(\mathbf{x}) \mid 将 \mathbf{y} \\ & \text{son} \\ & \text{str}(\mathbf{x}) \mid \mathcal{H} \\ & \mathbf{y} \\ & \text{str}(\mathbf{x}) \\ & \mathbf{y} \\$$

#### 算术运算符

x + 5 | 加, 计算结果为7. x - 5 | 减, 计算结果为-3.

x \* 5 | 乘, 计算结果为10.

x / 5 | 除, Python 2.x版本的计算结果为0, Python 3.x版本的计算结果为0.4.

x \*\* 2 | 暴运算, 即 x2, 计算结果为4.

False None 0 " " () [] {} | False值. and | 等价于 "&",表示 "且". or | 等价于 "| ",表示 "或". not |表示 "非".

#### 子行

S = u""|定义Unicode字符串。
S = r""|定义服的字符串,避免字符串中的字符转义,在正则表达式中经常使用到。

s = "cookdata" | 定义值为 "cookdata" 的变量s. len(s) | 返回s的字符个数8.

s.lower() | 将字符串s中的字母全部转换为小写. s.upper() | 将字符串s中的字母全部转换为大写. s.capitalize() | 将字符串s中的首个字符转换为大

写,其余部分转换为小写. s.replace('k', 'l') | 使用字符 "1" 替换掉s中所有的字符 "k", 返回结果为cooldata.

相明子传 K , 足面前水分COOlatta, s.s.trip() 土 全陸排棄衛 西町中最后面的空格. s.split("\t") | 使用刺表符 "\t" 分割字符率s. '%s is No.%d' %(s, 1) | 取出:的值和数值:成次放 入字符率%s is No.%d的相应位置, 返回结果为cookdata

'{} is No.{}'.format(s, 1) | 取出s的值和数值 1依次放入()相应位置,返回结果为cookdata is No.1.

#### 列夷

l = ['c', 'o', 'o', 'k', 1] | 创建一个包含字符元素c, o, o, k和整数1的列表.

list() | 创建空列表,或将其他数据结构转换为列表.

1[0] | 返回列表的第一个元素,即字符c. 1[-1] | 返回列表的最后一个元素,即1.

1[1:3] | 列表切片, 返回包含原列表的第二个元素和 第三个元素的列表['o','o'].

len(1) | 返回列表的元素个数. l[::-1] | 将列表进行逆序排列.

l.reverse() | 将列表进行进序排列。

l.insert(1, 'b') | 在指定的索引位置插入 'b' . l.append() | 在列表末尾添加元素.

1.extend(L) | 等价于"1+L",将列表L中的元素依次添加到10末尾。

1.remove() | 删除列表中的某个元素.

1.pop() | 等价于 "del 1[]", 删除列表中对应索引 位置的元素。

"..join(['c', 'o', 'o', 'k']) | 将列表中的 各个字符串元素用空格连接起来并转换为字符串,返回 结果为c o o k.

#### 又仟读与

f = open(filename,mode) | 返回一个文件对象f. 读文件 "model = r" , 写文件 "model = w" . f.read(size) | 返回包含前size个字符的字符串。 f.readline() | 每次读集一行,返回该行字符串。 f.readlines() | 返回包含整个文件内容的列表, 列 表的元素为文件的每一行内客所构成的字符串。 f.close() | 美国文件并释放它所占用的系统答案

with open("cookdata.txt", "r") as f: content = f.readlines() 在with主体块语句执行完后,自动关闭文件并释放占用 的系统管源。

import csv f = open("cookdata.csv", "r") csvreader = csv.reader(f) content\_list = list(csvreader) 该取csv文件,并配数程序接为一嵌套列表(列来的元 素仍是一个列表)content\_list.

#### 写入文件

读取文件

f.write(s) print(s, file=f) 两种等价的方式,将字符串s写入文件对象f中.

#### 1711

def sum(a, b=1): return a+b

命名为ob i, methodname的函数,

定义求和函数sum(),该函数要求输入位置参数a,带款 认值的参数b为可违参数,其默认值为1,函数返回结果 为a+b的计算结果。

sum(1, b=10) | 故行sum() 函數、 适回结果为11. def sum(\*args, \*\*kwargs) | 不是社多數、\*args 接收包含多个位置多数的元法。\*\*kwargs接收包含多个 类徵字参数的字典。 obj.methodname | 一个方法是一个"属于"对象并被

#### map() 和lamb

map(func, sequence) | 将函数依次作用在序列的每个元素上,把结果作为一个新的序列返回。 lambda a, b:a+b | 腰名函数,正常函数定义的语法 槽,a和b为输入参数,a+b为函数主体和返回的值。

#### 模均

import module as alias | 导入模块,并取一个射 名,使用alias. fune即可调用模块内的函数, from module import = | 导入模块中的所有函数, 直接使用函数名即可调用模块内函数。 from module import funcl, func2 | 导入模块 中的部分函数,使用funcl可直接调用设函数。

#### 面向对象的类

class Athlete(object):
 def \_\_init\_\_(self, name, age):
 self.name = name
 self.name = name
 self.name = name
 self.name = self.name(self):
 return self.name.capitalize()
 使用美鲜平class定义Athlete类,该类继承的ject类,
 物始充变量为self.name和self.age、并定义类的capitalize name为法带self.name参享并要发出可

a.name | 延囲aöname 県性、延囲 'james'.
a.age | 遅囲a参nage 黒性、延囲23.
a.capitalize\_name() | 項用a鲂实例方法capitalize\_name. 延囲 'james'.
isinstance(a, Athlete) | 判断a是否是Athlete类

a = Athlete('james', '23')|创建实例a.

#### 编码和解码

的实例.

ASCII | 基于拉丁字母的一套电脑编码系统,不包含中文、日文等非英语字符.

GBK | GBK兼容ASCII, 同时收录中文、日文等, 使用两个字节编码一个汉字.

Unicode | 牧录超过十万个字符,统一了所有语言文字的标准编码集,包括UTF-8和UTF-16两种实现方式.

S = u'中文'|定义Unicode字符串 '中文' S.encode('utf-8')| 使用UTF-8編碼集将Unicode字符串 s編码为str字符串。

S = '中文' s.decode('utf-8') 定义str字符串'中文',并解码为Unicode字符串.

import chardet chardet.detect(s) 检测字符串的编码方式.

### 异常处理

语法错误 | Syntax Errors, 代码编译时检测到的错误. 异常 | Exceptions, 代码运行时检测到的错误, 如类型

异常。Exceptions,代码近行时检测到的错误。如实型错误(TypeError)、数值错误(ValueError)、索引错误(IndexError)和属性错误(AttributeError)等。

#### try: statement except:

pOSS 先尝试运行try部分的statement语句,如果能够正常运 行,则跳过except部分,如果运行出现错误,则跳过错 误代码运行except部分的pass语句,"放过"错误.

#### try:

statement except Exception, e: print "Error Happened: %s" %e 指定要处理的运行时错误英型, 如果指定的错误出现。 则打印极错信息,e是发生错误时的具体错误信息.

#### t

statement1 except (Exception1, Exception2), e: statement2 #发生指定错误时的处理 else: statement3 #正确时运行

finally: statement4 #无论对错都运行 trv··except···else···finally句式。

### 抛出异常

raise Exception('0ops!') | 主动抛出异常Exception, 错误提示信息为'0ops!'.
assert statement, e | 若statement语句运行结果

assert statement, e | 若statement语句运行结果为True,则继续运行代码,否则拠出e的错误提示信息.

### 正则表达式

### 正则表达式模块 re

import re

raw\_s = r'\d{17}[\d\x]|\d{15}'
pattern = re.compile(raw\_s)
re.search(pattern, s)
用于匹配身份证号。

首先使用原始字符串定义正则表达式模式;然后编译原 始字符串为正则表达式Pattern对象;最后对整个字符串 s进行模式搜索,如果模式匹配,则近回MatchObject的 实例,如果该字符串没有模式匹配,则近回Mone.

re.search(r'\d{17}[\d\x]\d{15}', s) | 将 Pattern编译过程与搜索过程合而为一. re.match(pattern, s) | 从字符串s的起始位置匹配一个模式, 如果起始位置匹配不成功, 则眨回None.

re.findall(pattern, s) | 返回一个包含所有满足 匹配模式的子事的列表. re.sub(pattern, repl, s) | 使用替换字符串repl 替换匹配到的子字符串.

re.split(pattern, s) | 利用满足匹配模式的子串 将字符串s分隔开,并返回一个列表.

#### 日期处理

from datetime import datetime format = "%Y-%m-%d %H:%M:%S" | 審定日期卷式, 如 "2017-10-01 13:40:00". date\_s = datetime.strptime(s,format) 符日期字符串按照指定日期格式转换为datetime类.

将日期字符串接匯指定日期格式转换为datetime date\_s.year | 获取日期字符串中的年份。 date\_s.month | 获取日期字符串中的月份。 datetime.now() | 获取现在的日期和时间。

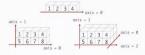




NumPy包是Python科学计算中的核心工具包之 一,它提供了高性能多维数组结构-ndarray和 用来操控这些数组的各种工具和函数.

导入包 (NumPy 1.14.0) import numpy as np

NumPy中的维度方向用axis表示.



#### 创建数组

np.zeros((2)) 2个0的一维数组. np.zeros((2, 3)) | 6个0的二维数组. np.ones((2, 3, 4)) | 24个1的三维数组. np.empty((3)) 创建未初始化的一维数组 np.eye(3) | 创建3×3的单位矩阵. np.array([[1,2], [3,4]]) | 二维数组. np.array([1,0], dtype=np.bool) | -维布尔数组.

np.array([[[1,2], [3,4]], [[5,6], [7,8]]]) | 三维数组.

np.full((2,3),7) | 数值全部为7的2乘3二 维数组.

np.aranae(4) 从0到3的一维数组, np.arange(5, 25, 5) 从5(包含)到 25 (不包含), 步幅为5的一维数组.

np.linspace(0, 9, 10) | 从0(包含)到 9(包含),总长为10的一维数组。

np.random.random((2, 3)) | 从区间[0.0, 1.0) 中随机抽取6个数,填入2乘3的数组中.

一个成为数据科学家的理由 www.hackdata.cn

np.cos(np.linspace(0, 2\*np.pi, 10)) 包含十个0到2π余弦函数值的一维数组.

#### 数组的结构

arr.shape 数组的维度. len(arr) 数组的长度. arr.ndim 数组维度的大小. arr.size | 数组中元素的个数. arr.dtype 数组数据的类型. arr.dtype.name 数组数据类型的名称. arr.astype('float') 将数据类型转换为 其它类型.

#### 数据类型

np.int64 | 64位整数型, np.float32 | 标准双精度浮点型. np.complex 由128位浮点代表的复数。 np.bool 布尔型. np.object | Python中object类型. np.nan NA值, 为浮点型 np.string\_ 定长字符串类型. np.unicode\_ | 定长unicode类型.

### 文件读取 np.save('my\_arr', arr) 以NumPy格式将

数组保存到本地, np.savetxt('my\_arr.txt', arr) | 以文 本格式将数组到本地, np.load('my\_arr.npy') | 读取本地NumPy 数组文件. np.loadtxt('my\_arr.txt') | 读取本地文 本数组文件, np.genfromtxt('file.csv', delimiter=',') 将其它文件文件转换成 NumPy数组.

#### 数组切片

说明: arr = np. arange(0,100). reshape(20,5)

arr[4] 提取第5个(行)的元素. arr[1,0] | 提取第2行第1列的元素. arr[[0,2], [1,3]] | 提取位于第一行第二 列[0.1]和第三行第四列[2.3]的两个元素。 arr[1][0] | 提取第2行第1列的元素. arr[0:5] | 提取第1行到第5行元素. arr[:5] | 提取前5行的元素. arr[0:10:2] | 提取第1, 3, 5, 7, 9行元素. arr[::-1] | 反序数组. arr[0:2, 2] 提取前2行的第3列的元素。 arr[:,:] | 提取所有行和所有列.

### 数组的操作

arr[4, ...] | 提取第5行的元素.

arr[arr<50] 提取所有小于50的元素.

说明: arr = np.array([[1,2],[3,4]])

arr.T 数组转置, [[1, 3], [2, 4]]. arr.ravel() 扁平化数组, [1, 2, 3, 4]. arr.reshape(4, 1) 重塑数组, 元素个数 需保持一致, [[1], [2], [3], [4]], arr.resize((2, 3)) | 重塑原数组, 元素个 数可不同, [[1, 2, 3], [4, 0, 0]]. np.append(arr, [[5,6]]) 扁平化后插 入新数组[[5,6]], 若指定axis参数, 则会按指 定方向插入, [1, 2, 3, 4, 5, 6]. np.insert(arr, 3, [5,6]) | 扁平化后 按索引插入新数组[[5,6]], 若指定axis参数, 则会按指定方插入, [1, 2, 3, 5, 6, 4]. np.delete(arr, 1) | 扁平化后按索引删除 元素,若指定axis参数则会按维度方删除,[1. 3, 4]. np.concatenate((arr, arr), axis=1)

按指定维度方向(需已存在)合并多个数组,

[[1, 2, 1, 2], [3, 4, 3, 4]].

np.hstack((arr, arr)) | 按水平方向合并 多个数组. np.vstack((arr, arr)) | 按竖直方向合并

多个数组.

np.split(arr,2) 按第0维度方向将数组切 成2份, [[1, 2]]和[[3, 4]].

np.split(arr,2,1) | 按第1维度方向将数组 切成2份, [[1],[3]]和[[2],[4]]. arr.tolist() | 将数组转换成列表.

#### 数组的运算

np.nan is np.nan | 返回True. np.nan == np.nan | 返回False. a == b | 逐元素比较,返回等维度布尔数组. a < 2 | 逐元素比较,返回等维度布尔数组. a+b, a-b, a\*b, a/b | 逐元素进行加、减、 乘、除运算。

a.dot(b) 数量积运算. a.max(a), a.min(a), a.mean(a), np.std(a) 求数组整体的最大、最小、均值 和标准差.

a.max(axis=0), a.min(axis=1) | 按指定 维度方向求最大值、最小值.

a.sort(axis=0) 按指定维度方向增序排序 数组.

a.copy() 深拷贝数组. a.view() 创建数组的view.

np.copy(a) | 浅拷贝数组.

np.array\_equal(a, a) | 数组间整体比较, 返回单个布尔值.

np.exp(a), np.sqrt(a), np.sin(a), np.log(a) | 逐元素地做幂、开平方、正弦、 自然对数运算.

在NumPv 1.14版本中不支持的运算包括: a.median()和a.corrcoef().





## り数据啼客

基于NumPy构建,利用它的高级数据结构和操作工具,可使数据分析工 作变得更加便捷高效.

#### 符号标记

S | 一个Series对象. df | 一个DataFrame对象。 导入包 (Pandas 0.20.1) import numpy as np import pandas as pd

#### 基本数据结构

#### Series

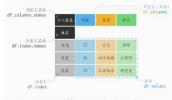
一维数据结构,包含行索引和数据两个部分.



s = pd.Series([14, 15, 17], index=[u'张某', u'李某', u'段某'])

#### DataFrame

二维数据结构。包含带索引的多列数据、各列的数据类型可能不同。



df = pd.DataFrame([[22, u'北京', u'律师'], [26, u'四川或都', u'工程师'], [24, u'江苏南京', u'研究员']], index=[u'张某', u'孝某', u'段某'], columns=[u'年龄', u'精贯', u'职业'])

从文件中读取数据 (DataFrame)

pd.read\_csv() | 从CSV文件读取. 参数: file, 文件再任: sep, 分陽符; index\_col, 用件行索引的一列或者 多列: usecols, 世界列: encoding, 字符编码类型, 通常为'utf-8'. pd.read\_table() | 从制表符分隔文件读取,如TSV. pd.read\_excel() | 从Excel文件读取 pd.read sal() 从SQL表或数据库读取 pd.read\_json() | 从JSON格式的URL或文件读取.

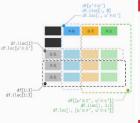
### 将DataFrame写入文件

df.to\_csv() 写入CSV文件 参数: file, 文件等册; sep, 分隔符; columns, 可入文件的列; header, 是否可入表头; index, 是否可入索引,

df.to\_excel() | 写入Excel文件, df.to\_sql() 写入SQL表或数据库. df.to\_json() 写入JS0N格式的文件。 df.to\_clipboard() 写入剪切板。

pd.read\_clipboard() | 从剪切板读取.

#### 数据索引



df[5:10] | 通过切片方式选取多行. df[col\_label] or df.col\_label | 选取列. df.loc[row\_label, col\_label] | 通过标签选取 df.sort 行/列. df.iloc[row\_loc, col\_loc] | 通过位置(自然 df.sort 数)选取行/列。

df.info

df[u'

业"是否

df[df[u df.loc[数据融合 (merge)

表达式等

df.filt

df.drop

df.drop

df.drop

df.sele

clude) 司

inplace=

axis=0, i

df.astv

df.inde

df.colu

df.set\_

df.rese

df.rein

据index和

df.rena

索引或列

df.sort

进行升序:

据"A" 及 df.rank

s.unique() 唯一值. s.value\_counts() 唯一值及其计数. df.head(n) | 前n行数据.

数据合并

张某 北京

8 中華 四月成都

9 我基 江东南京

82 48 BE

1 \$2 M0 wright

2 RA 240 NaN

HZ 40 H1

1 0-E 200 m ( A to

pd.merge(left, right) | 类数据库的数据融合操作.

参数: how, 融合方式, 包括左连接、右连接、内连接(聚认) 和外连接; on, 连

推帧: left\_on, 左锁: right\_on, 右键: left\_index, 是否将left行索引作

为左键; right\_index, 是否符right行索引作为右键.

母某 22

李某 26

2 EE 24

HA YM 格質

1 42 2 HAAR

df,tail(n) | 后n行数据, df.shap



数据啼客

数据融合 (join)

张某 北京 李某 10月或都 我基 江苏南京

8 2 22 EX 李基 26 四月成年 RE 24 NAN

FE 20

left.join(right) 在索引上的数据融合操作.

#### 数据融合 (combine first)

left.combine\_first(right) 在数据融合的同时 (行索引和列索引取并集), 使用right 中的值填补left中相应位置的缺失值.

#### 轴向连接





pd.concat([df1, df2]) 轴向连接多个DataFrame.

### 轴向旋转和数据转换

轴向旋转

df.stack() 将数据的列"旋转"为行。 df.unstack() | 将数据的行"旋转"为列.

#### 数据转换

回一个GroupBy对象,

s.map() 利用函数或映射进行数据转换, df.applymap() 利用函数或映射进行数据转换. df.replace() | 替换元素值. df.columns.map(str.upper) | 将列名转换为大写.

### 数据聚合与分组运算 df.groupby('A') 根据 "A" 列的值进行分组,并返

df.groupby(['A', 'B']) | 根据 "A" 列和 "B" 列的 值组合进行分组,并返回一个GroupBy对象. df.groupby('A')['B'].mean() | 根据 "A" 列的值 进行分组,并计算每一组"B"的均值, df.groupby('A').agg(np.mean) 根据 "A" 列的 值进行分组,并计算每一组中各列的均值. df.apply(np.mean) | 对DataFrame的每一列求均值. df.apply(np.mean, axis=1) | 对DataFrame的每一 行求均值. df.pivot\_table(index='A', values=['B', 'C'] , agafunc=mean) 创建数据透视表,根据"A"列的

#### 离散化和喷变量编码

值进行分组,并计算每一组中"B"和"C"的均值.

pd.cut(s, bins=[0, 5, 10], labels=['A', 'B']) 离散化和面元划分. pd.qcut(s, 2, labels=['A', 'B']) | 等頻离散化. pd.get\_dummies(df, columns=['A'],

drop first=True) 对"A"列进行哑变量编码,并去掉第一个编码特征。

### 串联方法

通过调用多个Pandas方法将多个处理过程串联起来 pd.merge(left, right, how='outer', indicator=True)

.query('\_merge == "left\_only"') .drop(['\_merge'], axis=1) 保留在left中但不在right中的行。

#### 文本数据规整

s.str.method

矢量化字符串方法,可跳过缺失值, 週缺失值将不报错

s.str.lower() 将所有元素转换为全部小写。

s.str.upper() 将所有元素转换为全部大写。 s.str.replace() | 替换值.

s.str.split(" ").str[0] | 空格分割, 并取列表中 的第一个元素,

s.str.split("").str.get(0) | 空格分割,并取列 表中的第一个元素,

s.str.split(" ", expand=True) | 空格分割,并 扩展为多列,

s.str.contains(r'\d') | 利用正则表达式判断是否 元素包含数字

s.str.extract(r'(\d.)', expand=False) | 利用

正则表达式提取每个元素中的数字。

#### 数据可视化

df.plot(kind='hist') 设置kind参数绘制直方图. df.plot.line() 调用对象接口绘制折线图. df.plot.scatter('x', 'y') | 绘制x-y散点图. 包括图形: W 机 Edf.plot.area: 从向柱状 Edf.plot.bar: 横向柱状 图df.plot.barh: 全图df.plot.box; 种核密度图df.plot.kde; 蜂具 Mdf.plot.hexbin: # Mdf.plot.pie.

### 一个成为数据科学家的理由

## Scikit-learn







## 数据 赔



Scikit-learn (0.19.1) 是基于NumPy、SciPy和Matplotlib的开源 Python机器学习包,它封装了一系列数据预处理、机器学习算法、模型 选择等工具,是数据分析师首选的机器学习工具包。

符号标记 X\_train | 训练数据. X\_test | 测试数据. X | 完整数据.

y\_train | 训练集标签. y\_test | 测试集标签. y | 数据标签.

#### 基本建模流程

②号入工具を from sklearn import datasets, preprocessing from sklearn.model\_selection import train\_test\_split from sklearn.lineer\_model import tinearRegression from sklearn.metrics import rZ\_score か 物料を

boston = datasets.load\_boston()

X = boston.data
y = boston.target

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.3)

① 数据标题

scaler = preprocessing.StandardScaler().fit(X\_train)

X\_train = scaler.transform(X\_train
X\_test = scaler.transform(X\_test)

X\_test = scaler.transform の 権型和非当組合

lr = LinearRegressio

lr.fit(X\_train, y\_tra ® 模型預測与评价

y\_pred = lr.predict(X\_test)
r2 score(y test, y pred)

#### 加载数据

✓ Scikit-learn支持以NumPy的arrays对象、Pandas对象、SciPy的稀疏 矩阵及其他可转换为数值型arrays的数据结构作为其输入。前提是数据 必须是数值型的。

✓ sklearn.datasets模块提供了一系列加载和获取著名数据集加商尾花、波士顿房价、Olivetti人脸、MNIST数据集等的工具,也包括了一些toy data如S型数据等的生成工具.

from sklearn.datasets import load\_iris

iris = load\_iris()

X = iris.data

y = iris.target

#### 训练集-测试集划分

stratify=y, test\_size=0.3) 特完整数据集的70%作为训练集.30%作为测试集.并使得测试集和训练 集中各类别数据的比例与原始数据集比例一级(stratify分层策略). 另外可通过设置 shuffle=True 提前打乱数据.

#### 数据预处理

#### 标准化

from sklearn.preprocessing import StandardScaler ① 构建榜换器实例

scaler = StandardScaler()

② 拟合及转换

scaler.fit\_transform(X\_train)

部分数据预处理方法	对应的sklearn的类
最小最大标准化	MinMaxScaler
One-Hot编码	OneHotEncoder
归一化	Normalizer
二值化(单个特征转换)	Binarizer
标签编码	LabelEncoder
缺失值填补	Imputer
多项式特征生成	PolynomialFeatures

#### 特征选择

from sklearn import feature\_selection as fs fs.SelectKBest(score\_func, k) | 过滤式 (Filter), 保留 得分接名前k的特征 (top k方式).

#### 有监督学习算法

from sklearn.linear\_model import LinearRegression ① 柏建模型字例

lr = LinearRegression(normalize=True)

② 训练模型 lr.fit(X\_train, y\_train)

3 作出預測 y\_pred = lr.predict(X\_test)

部分回归算法	对应的sklearn的类
LASS0	linear_model.Lasso
Ridge	linear_model.Ridge
ElasticNet	linear_model.ElasticNet
回归村	tree.DecisionTreeRegressor

#### 45.89

from sklearn.tree import DecisionTreeClassifier
clf = DecisionTreeClassifier(max\_depth=5)
clf.fit(X\_train, y\_train)

y\_pred = clf.predict(X\_test)
y\_prob = clf.predict\_proba(X\_test)

决策树分类算法示例,对于二分类问题,y\_prob为每个样本预测 为"0"和"1"类的概率。

部分分类算法	对应的sklearn的类
逻辑回归	linear_model.LogisticRegression
支持向量机	svm.SVC
朴素贝叶斯	naive_bayes.GaussianNB
K近邻	neighbors.NearestNeighbors

#### 集成

sklearn.ensemble模块包含了一系列基于集成思想的分类、回归和 离群值检测方法。

from sklearn.ensemble import RandomForestClassifier clf = RandomForestClassifier(n\_estimators=20) clf.fit(X\_train, y\_train)

y\_pred = clf.predict(X\_test)
y\_prob = clf.predict\_proba(X\_test)

y\_prob = clf.predict\_proba(X\_test) 部分集成算法 对应的sklearn的类 AdaBoost ensemble.AdaBoostClgssifier

ensemble.AdaBoostRegressor 基于梯度提升 ensemble.GradientBoostingClassifier ensemble.GradientBoostingRegressor

#### 无监督学习算法

#### 聚类

sklearn.cluster模块包含了一系列无监督聚类算法.

from sklearn.cluster import KMeans ① 构建聚类实例

kmeans = KMeans(n\_clusters=3, random\_state=0) ② 拟合

kmeans.fit(X\_train)

⑤ 预测

kmeans.predict(X\_test)

部分聚类算法	对应的sklearn的类
DBSCAN	cluster.DBSCAN
层次聚类	cluster.AgglomerativeClustering
谱聚类	cluster.SpectralClustering
谱聚类	cluster.SpectralClustering

#### 降维

from sklearn.decomposition import PCA
pca = PCA(n\_components=2)
pca.fit\_transform(X\_train)

通过主成分分析 (PCA) 将原始数据映射到2维空间,

部分降维方算法	对应的sklearn的类
线性判别分析	discriminant_analysis
	LinearDiscriminantAnalysis
核主成分分析	decomposition.KernelPCA
局部线性映射	manifold.LocallyLinearEmbeddin
t-SNE	manifold.TSNE

#### 模型评价

sklearn.metrics模块包含了一系列用于评价模型的评分函数、损失函数以及成对数据的距离度量函数.

from sklearn.metrics import accuracy\_score metrics.accuracy\_score(y\_true, y\_pred) 对于測试集而言,y\_test即是y\_true,大部分函数都必須整合真实值y\_true和預測

值y\_pred, 基于排放的关系,以下不再重复写这两个参数。

#### 回归模型评价

metrics.mean\_absolute\_error() | 平均絶対误差MAE. metrics.mean\_squared\_error() | 均方误差MSE. metrics.r2\_score() | 決定系數R<sup>2</sup>.

#### 分类模型评价

metrics.accuracy\_score() | 正確率. metrics.precision\_score() | 本美精磯平, metrics.fl\_score() | Fi组. metrics.log\_loss() | 対数很失成文叉薄損失. metrics.confusion\_metrix | 温声矩序. metrics.cossification\_report | 本参种评价的分素提売.

#### 交叉验证及超参数调优

#### 交叉规则

使用5折交叉验证对决策树模型进行评估,使用的评分函数为F1值.  $\checkmark$  sklearn提供了部分帶交叉验证功能的模型类如LassoCV、LogisticRegressionCV等,这些类包含cv参数.

#### 超参数调优---网格搜索

from sklearm.model\_selection import GridSearch(V svc = svm.SVC): ['linear', 'rbf'], 'C':[1, 10]] grid\_search = fridSearch(V(svc, poroms, cv-5) grid\_search.fit(X\_train, y\_train) grid\_search.best\_poroms. 在参数网络上世行穷年授末、方法简单但是搜索进度慢(超季数较多 时),且不愿多数较多数发明中的局部都是。

#### 超参数调优——随机搜索

from sklearn.model\_selection import RandomizedSearchCV from scipy.stats import randint svc = svm.SVC()

param\_dist = {'kernel':['linear', 'rbf'], 'C':randint(1, 20)}
random\_search = RandomizedSearchCV(svc, param\_dist, n\_iter=10)
random\_search | host params

在参数子空间中进行随机搜索,选取空间中的100个点进行建模(可从 scipy,stats常见分布如正恋分布norm,均匀分布uniform中随机采样 得到),时间耗费较少,更容易找到局部最优。





## Matplotlib



plt.pie()

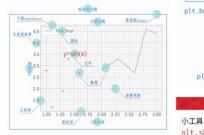




Matplotlib是Python中一个非常有名的绘图包, 它能快速绘制具有印刷品质的图形, 并且支持跨 平台运行和交互式操作.

导入包 (matplotlib 2.1.1) import matplotlib.pyplot as plt

图形的各元素名称如下:



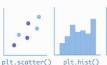
#### 基本概念

绘图框是图形的最高容器, 所有图形必须放置在 绘图框中

子图 是绘图框中所包含的图形,即便绘图框只 包含一幅图,也称之为子图.

元素 是组成子图的部件,从子图最内部的数据 线条到外围的坐标轴标签等都属于元素.

#### 图形样式





文字调整

8的绘图框.

刻度调整

plt.bar()

plt.barh()

plt.boxplot() plt.fill() plt.hexbin()

plt.violinplot() plt.contour()

图形设置

plt.savefig('fig.eps') | 保存图形.

plt.xlim(0,1) 设置x轴显示范围, y轴同理.

plt.tick\_params(axis='x', size=50,

labelsize=20) | 修改主刻度及其标签样式.

plt.grid(which='major') 添加背景网格.

plt.xscale('log') 设置轴比例, y轴同理.

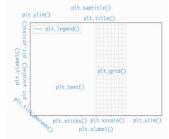
plt.xticks(np.arange(0,1,0.2)) 修改

plt.show() 显示图形.

主刻度标签范围, y轴同理.

plt.title('title') | 添加子图标题. plt.legend(['ln','pt']) | 添加图例. plt.xlabel('x') | 添加x轴标签, y轴同理. plt.text(1,1,'t') 在(1,1)位置添加文字.

plt.annotate('t', xy=(1,2), xytext =(3,4), arrowprops=dict()) | 在子图中 添加带有指向箭头的文字注释. plt.suptitle('T') |添加绘图框标题.



from matplotlib.font\_manager import FontProperties myfont=FontProperties(fname='FangSong.TTF') plt.title(u'简单的正弦图形', fontpro-

#### perties=myfont) 使用中文字体 plt.figure(figsize=(12, 8)) | 创建12×

#### 绘制子图

plt.subplot(1,2,1) | 创建1 × 2的子图矩 阵, 当前绘制第一幅子图.

高级绘图

plt.tight\_layout() | 自动调整子图问距. plt.subplots\_adjust() | 调整子图大小.

#### 高级操作

fig = plt.figure() | 创建绘图框并返回绘 图框对象。

ax = fig.add\_subplot(121) | 对创建好 的绘图框对象添加1×2的子图矩阵并返回第一个 子图对象,

fig, ax = plt.subplots(1,2) 创建 $1 \times 2$ 的

子图矩阵并同时返回绘图框和两个子图对象, 其 效果等于前两者之和.

ax.plot(), ax.bar(), ... 向子图中添 加图形, 其余样式同理,

ax.set\_title(), ax.set\_xlabel(),... 对子图添加标题、标签等,其余元素同理,

plt.gca() 返回当前绘图框中的子图对象. plt.gcf() 返回当前绘图框对象.

ax.invert\_xaxis() 反序x轴, y轴同理.

ax.xaxis.tick\_top() | 将x轴主刻度标签移 动到上方, y轴同理.

ax.tick\_params() 设置坐标轴刻度.

ax.set\_frame\_on(False) | 关闭子图边框. ax.set\_axis\_off() | 关闭子图所有坐标轴.

from mpl\_toolkits.mplot3d import Axes3D

fig = plt.figure()

ax = Axes3D(fia) | 创建3D子图对象。

### 常用颜色

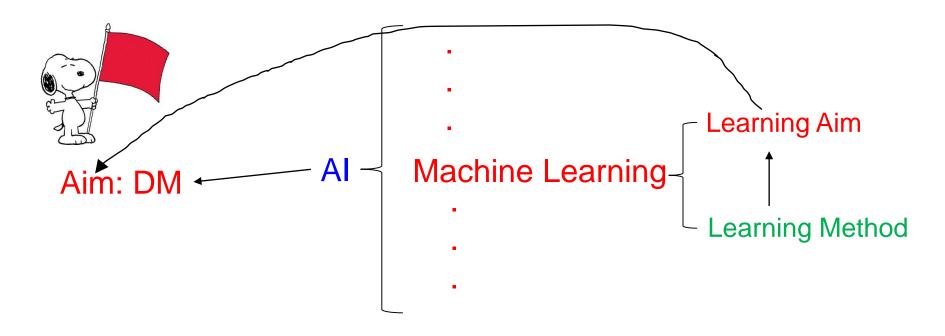




## Review



## How do you understand DM、AI、ML?



## Review



# **Data Mining Methods**

### ■ 分类

- 分类问题属于预测性的问题,它跟普通 预测问题的区别在于其预测的结果是类 别。
- 典型三种算法: Logistic回归; 决策树; 神经网络
- 应用场所: 判定类别等

### ■ 聚类

- 指把一组个体按照相似性归成若干类别。
- 典型算法: 基于欧氏距离; K-means算法
- 应用场所: 市场细分

### ■ 关联

- 挖掘发现大量数据中项集之间有趣的关 联或相关联系
- 典型算法: Apriori算法; FP-Growth算法:
- 应用场所:交叉销售(购物篮-啤酒与尿布)



# **DATA ANALYTICS:**

# DATA MINING AND BIG DATA



— Mining Rules 1



## **Association Rules**







Frequent Item\_sets 频繁项集

Support 支持度

Confidence 置信度

Support  $(X => Y) = |\{X, Y\}| / |\text{Items}| \text{ Confidence } (X => Y) = |\{X, Y\}| / |X_{\text{Items}}|$ 



In 10 days, there were 4 days when Xiao Ming ate a hamburger and bought a watermelon at the same time. What is support? ( )

A 40%

B 20%

50%

D 60%

In 10 days, Xiao Ming ate hamburgers for 8 days. During the 8 days, Xiao Ming also bought watermelon for 4 days. What is the confidence? ( )

A 40%

B 20%

50%

**60%** 

## **Association Rules**



Support: Find Frequent Item\_sets

Confidence: Find Association rules

## **Association Rules**

篮子	育品1	育品2	育品3
1	香草威化	香蕉	狗粮
2	香蕉	面包	酸奶
3	香蕉	苹果	酸奶
4	香草威化	香蕉	生奶油
5	面包	香草威化	酸奶
6	牛奶	面包	香蕉
7	香草威化	苹果	香蕉
8	酸奶	苹果	香草威化
9	香草威化	香蕉	牛奶
40	4 th	<b>ア</b> /2	at at 30



Confidence (Vanilla Wafer => Banana) = 4/6 = 67%

Confidence (Banana => Vanilla Wafer) = 4/8 = 50%

# Apriori算法



First, Set a Support initial threshold.

Set 1- Dataitems list.

Set 2- Dataitems list.

Set 3- Dataitems list.

. . . . . .

Set n- Dataitems list.

1994 — "Fast algorithms for mining association rules in large databases"

# Objectives (课程目标)



O1: Master Machine Learning classic algorithm

O2: Master classic Statics Methods

O3: Data Mining Methods

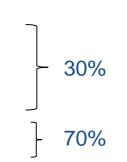
O4: Solve some Business examples

W1: Statics W2: Statics Quiz1

W3: Machine Learning W4: Machine Learning

W5: Machine Learning Quiz2 W6: Data Ming Rules

W7: Spark + Dig Data Quiz 3W8: Example Final Exam







贵在坚持!