

# 如何打一个数据挖掘比赛 (入门版)

## 贡献者名单

### 代码贡献

牧小熊、司玉鑫

### 代码测评

潘姝宇、王振东、骆秀韬

## 贡献组织名单

厦门大学WISERCLUB-竞赛部



这是一份简易的竞赛教程，我们的目的是帮助同学们迈出AI训练大师之路的第一步。数据挖掘中会有很多需要学习的地方，建议入门的同学可以暂时不用着急去弄各个代码的原理，先跑通代码，然后看代码中涉及的知识点去查询相关资料进行学习，这样能让你学习更加有目标性，也容易找到学习的乐趣。千里之行，始于足下，从这里，开启你的AI学习之旅吧！

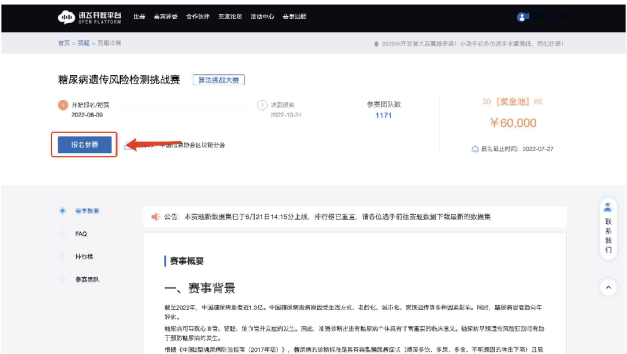
## 1.准备步骤

### 平台注册与比赛报名

- 1. 点击赛事链接：<https://challenge.xfyun.cn/topic/info?type=diabetes&ch=ds22-dw-wd06>
- 2. 注册（记得填写个人信息）



- 3. 点击报名参赛，显示成功报名



### 数据下载

- 1. 数据获取
  - 官网下载数据：[下载数据及实名验证](#)
  - 请把数据文件和代码文件放在同一个文件夹下，保证正常运行

## 环境配置

- python环境的搭建请参考
  - Mac设备：[Mac上安装Anaconda最全教程](#)
  - Windows设备：[Anaconda超详细安装教程\(Windows环境下\)](#)
  - [Anaconda的介绍与安装](#)，[jupyter notebook安装和使用](#)

## 2. 实践思路

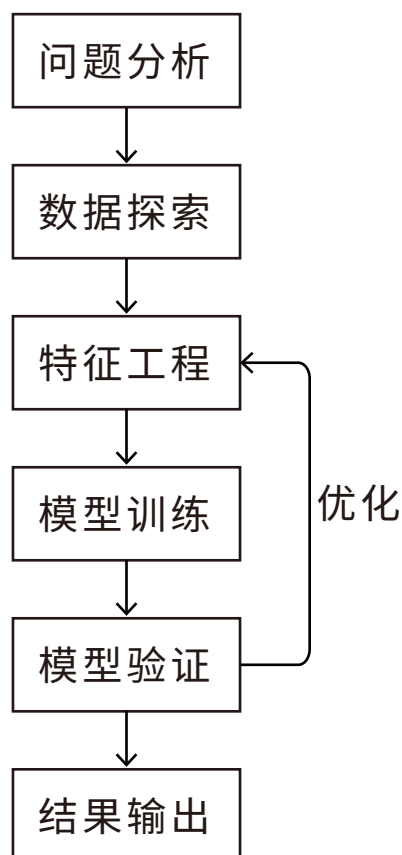
### 思路讲解

本次比赛是一个数据挖掘赛，需要选手通过训练集数据构建模型，然后对验证集数据进行预测，预测结果进行提交。

本题的任务是构建一种模型，该模型能够根据患者的测试数据来预测这个患者是否患有糖尿病。这种类型的任务是典型的二分类问题（患有糖尿病/不患有糖尿病），模型的预测输出为0或1（患有糖尿病：1，未患有糖尿病：0）。

机器学习中，关于分类任务我们一般会想到[逻辑回归](#)、[决策树](#)等算法，在这个Baseline中，我们尝试使用决策树来构建我们的模型。

我们在解决机器学习问题时，一般会遵循以下流程：



## 代码实现

```
1  #安装相关依赖库 如果是windows系统,cmd命令框中输入pip安装,参考上述环境配置
2  #!pip install sklearn
3  #!pip install pandas
4  #-----
5  #导入库
6  #-----数据探索-----
7  import pandas as pd
8  from sklearn.tree import DecisionTreeClassifier
9  #数据预处理
10 data1=pd.read_csv('比赛训练集.csv',encoding='gbk')
11 data2=pd.read_csv('比赛测试集.csv',encoding='gbk')
12 #label标记为-1
13 data2['患有糖尿病标识']=-1
14 #训练集和测试机合并
15 data=pd.concat([data1,data2],axis=0,ignore_index=True)
16 #将舒张压特征中的缺失值填充为-1
17 data['舒张压']=data['舒张压'].fillna(-1)
18
19 #-----特征工程-----
20 """
21 将出生年份换算成年龄
22 """
23 data['出生年份']=2022-data['出生年份'] #换成年龄
24
25
26 """
27 人体的成人体重指数正常值是在18.5-24之间
28 低于18.5是体重指数过轻
29 在24-27之间是体重超重
30 27以上考虑是肥胖
31 高于32了就是非常的肥胖。
32 """
33 def BMI(a):
34     if a<18.5:
35         return 0
36     elif 18.5<=a<=24:
37         return 1
38     elif 24<a<=27:
39         return 2
40     elif 27<a<=32:
41         return 3
42     else:
43         return 4
44
45 data['BMI']=data['体重指数'].apply(BMI)
```

```

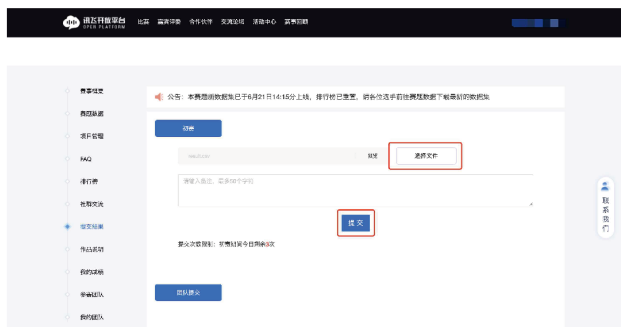
46
47 #糖尿病家族史
48 """
49 无记录
50 叔叔或者姑姑有一方患有糖尿病/叔叔或姑姑有一方患有糖尿病
51 父母有一方患有糖尿病
52 """
53 def FHOD(a):
54     if a=='无记录':
55         return 0
56     elif a=='叔叔或者姑姑有一方患有糖尿病' or a=='叔叔或姑姑有一方患有糖尿病':
57         return 1
58     else:
59         return 2
60
61
62 data['糖尿病家族史']=data['糖尿病家族史'].apply(FHOD)
63 """
64 舒张压范围为60-90
65 """
66 def DBP(a):
67     if a<60:
68         return 0
69     elif 60<=a<=90:
70         return 1
71     elif a>90:
72         return 2
73     else:
74         return a
75 data['DBP']=data['舒张压'].apply(DBP)
76
77
78 #-----
79 #将处理好的特征工程分为训练集和测试集,其中训练集是用来训练模型,测试集用来评估模型准确度
80 #其中编号和患者是否得糖尿病没有任何联系,属于无关特征予以删除
81 train=data[data['患有糖尿病标识'] !=-1]
82 test=data[data['患有糖尿病标识'] ==-1]
83 train_label=train['患有糖尿病标识']
84 train=train.drop(['编号','患有糖尿病标识'],axis=1)
85 test=test.drop(['编号','患有糖尿病标识'],axis=1)
86
87 #-----模型训练-----
88 model = DecisionTreeClassifier()
89 model.fit(train, train_label)
90 y_pre=model.predict(test)
91 y_pre
92
93 #-----结果输出-----
94 result=pd.read_csv('提交示例.csv')

```

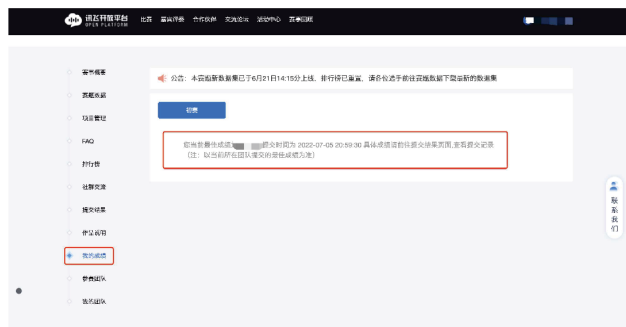
```
95 result['label']=y_pre
96 result.to_csv('result-de.csv',index=False)
```

## 结果提交

- 在提交结果处提交，提交 `预测结果.csv`，查看自己的成绩排名。



选择刚才生成的result.csv点击提交



点击我的成绩查看结果

## 3. 学习提升

数据挖掘流程及参考学习资料：[xunfei\\_demo.ipynb](https://github.com/xunfei-demo.ipynb)（用jupyter notebook打开）

## 4. 提问

如您遇到更多问题，可在Q&A文档中记录，我们会跟进反馈

<https://docs.qq.com/sheet/DUEVNBENIT1ZiZ2t2?tab=BB08J2>

## 评测贡献者

逸杰 为光清醒 清幻苗隼 Alex 小力 余舍 whl  
Joo若辰 Ww 夜浪 伊小雪 CNWL 驹驹 yong Clairezstar 凸 康书豪 CELFS Bib8o David MDR R007 江晓 刹那 alcm 奋斗 向辉  
北笙是光。 Jzkuan 乔森 kongla Carolyn crown 夜初晨 比诺莫 jchhhh eureka 李中平 果果红 朱子静 夏漫漫 秦小喵 柏杨  
戳戳戳 JRong. 只争朝夕 风 昭日月明 晴空万里 不畏 Amber 小左皮特 wltr Beyond 隔壁小王 诗仙李白 小包 章鱼情圣 Elohim  
Meimei 睡觉大王 风中有音 TTTT 谦 不忘初心 黄色枫叶 beckygong 隐隐约约 南栀倾寒 ZIZOU CX 白米饭 槐月初叁 Ch-watery 暮月矜心  
静水流深 四喜-VUW loong Angie#4370 JFDXP-中科院 像风一样自由 Abouerp Vellichor 薛定猫的谔 write-1 蜗一口吃鲨鱼 Fridayssss  
future 夷弥的藕饼 纵使明月照沟渠 路从今夜白 whitebloodcell 碧沼 William 王梓 onejoyaddition 凡凡 长安一片月 Never give up epoch  
。 益祥(Eason) yzpirate\_jackson BC. Meeerlin\_Wong SADAME-研一 不想学就完蛋。 Orient520520 不识食者 樂 好好学习天天向上 Crispo  
走成华大道的阿宝 Wilson Edwards 棒棒不是糖 看到我，请叫我去学习运动！ CheeseBoozhangting\_tbsi21 sapient\_river 哎哟哥哥的阳光宅男  
Enguane1 热爱分析的黎曼 栋第二十八年夏至 慈我的小名叫偶的 A.I StarryPilgrim 曹伟 w857736758 suffice 最幸运的幸福 Sundeyscoming  
静安其变 xg headlights 西红柿爱喝水 橙子 hellohaozheng 德艺双馨淡 要瘦瘦瘦呀 NaN 根本睡不着 仅自己可见 liyang4979 海珀克利特  
夏日回音 elephant ECCUSXR 友培 静君啊 waitzkin KevinDavis paletteboo 洛希极限 张璇 公大 say 嗨 try-except MYDLWZY Immelon TIAN  
iLQ0110 余一十五 lee 是非之欢 秋兰为佩 彩笔不彩 Ivy.X Monarch 热合帕提 JUSTHS 谦 白云千仞 忽而今夏 ZCLYiqun Liu 樱菲沐槿 拂浪  
李香兰 J.Y.Z.Z AI 小菜机 Snowya 都一曾小胖啊 Davidlp 陈佳莹 Laswell 翔 王jiajia vampire 卡卡南安 麒麟阁下 南波 onequiller7 Rory  
余泰 00cary 萝卜苗 薛喜悦 夏至 馨霏儿 jackin CYFeeYong 胖原体 陈亦余 youran 五斗米 DAMON Trank 普通人 咸鱼干 知不知道 长沙 小胖  
林言 无盐 J&J 夜殇 wm 知非 Eric 轻语 初见 士千 受戒 paul 靖言 许椰树\_zzzz EASY Devin 珞索 阿达 太和 张清 小胖  
mr2 yzq 小满 哇塞 祇辰 十六 abc 听瑜 Min

# 评测贡献高校

浙江大学 云南大学 天津科技大学 物资学院 西安石油大学 西安财经大学 临沂大学 中南民族大学 南方医科大学 江苏海洋大学 昆明理工大学 新疆理工学院 香港中文大学 南京财经大学 浙江财经大学 重庆邮电大学 燕山大学 太原理工大学 华北电力大学 河北农业大学 湖北工业大学 大连理工大学 济南大学 广州华立学院 湖南师范大学 华东师范大学 广东东软学院 华南农业大学 北京理工大学 广西师范大学 哈尔滨商业大学 新疆大学 哈尔滨理工大学 天津中医药大学 北京工业大学 美国埃默里大学 桂林理工大学 上海大学 安徽医科大学 Columbia University 中国科学技术大学 暨南大学 波士顿大学 河海大学 中山大学 上海应用技术大学 上海工程技术大学 成都信息工程大学 上海电力大学 东北大学 亚利桑那州立大学 桂林电子科技大学 东华大学 中国科学院大学 东北财经大学 山东财经大学 河北政法职业学院 中北大学 安徽师范大学 湖北大学 黑龙江大学 沈阳航空航天大学 同济大学 华南理工大学 西安交通大学 河北工业大学 广东医科大学 长安大学 杭州电子科技大学 长春大学 北京师范大学-香港浸会大学联合国际学院 东北林业大学 北京交通大学 北京科技大学 天津商业大学 郑州大学 北京邮电大学 青岛科技大学 华东理工大学 西安电子科技大学 中国科学技术信息研究所 中国科学院大学软件研究所 大连海事大学 广东工业大学 福州大学 上海交通大学 华北理工大学 哈尔滨工业大学 华中科技大学 上海财经大学 成都大学 吉林农业大学 复旦大学 中国石油大学华东 清华大学 江西理工大学 广州大学 沈阳理工大学 重庆理工大学 武汉邮电科学研究院 科罗拉多州立大学 北京石油化工学院 安徽工程大学 广西民族大学 中国人民公安大学 哈尔滨工程大学 中国石油大学(北京) 天津大学 重庆对外经贸学院 伦敦政治经济学院 湖南大学 北部湾大学 湖北警官学院 广东金融学院 大连交通大学 深圳大学 大连民族大学 东南大学 长沙理工大学 成都理工大学 圣安德鲁斯大学 河南工业大学 山西大学 长江大学 华东交通大学 湖南工商大学 河南城建学院 兰州财经大学 广州商学院 兰州交通大学 曼彻斯特大学 武汉理工大学 吉林大学 金陵科技学院 云南工商学院 杜伦大学 曲阜师范大学 维多利亚大学 浙江科技学院 浙江传媒学院 南京工业大学 齐鲁工业大学 三峡大学 武汉大学 中南大学 海南大学 莆田学院 天津财经 文华学院 中国矿业大学 南开大学 五邑大学 扬州大学