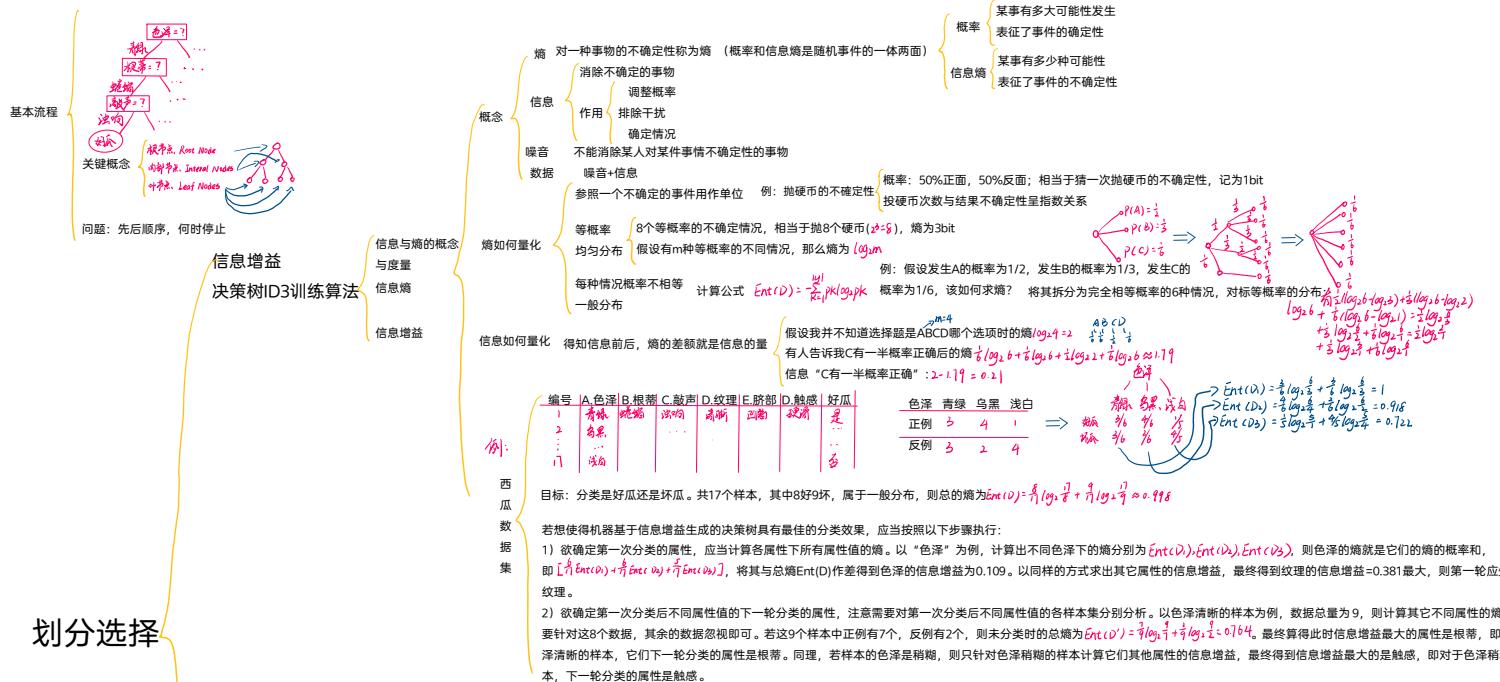
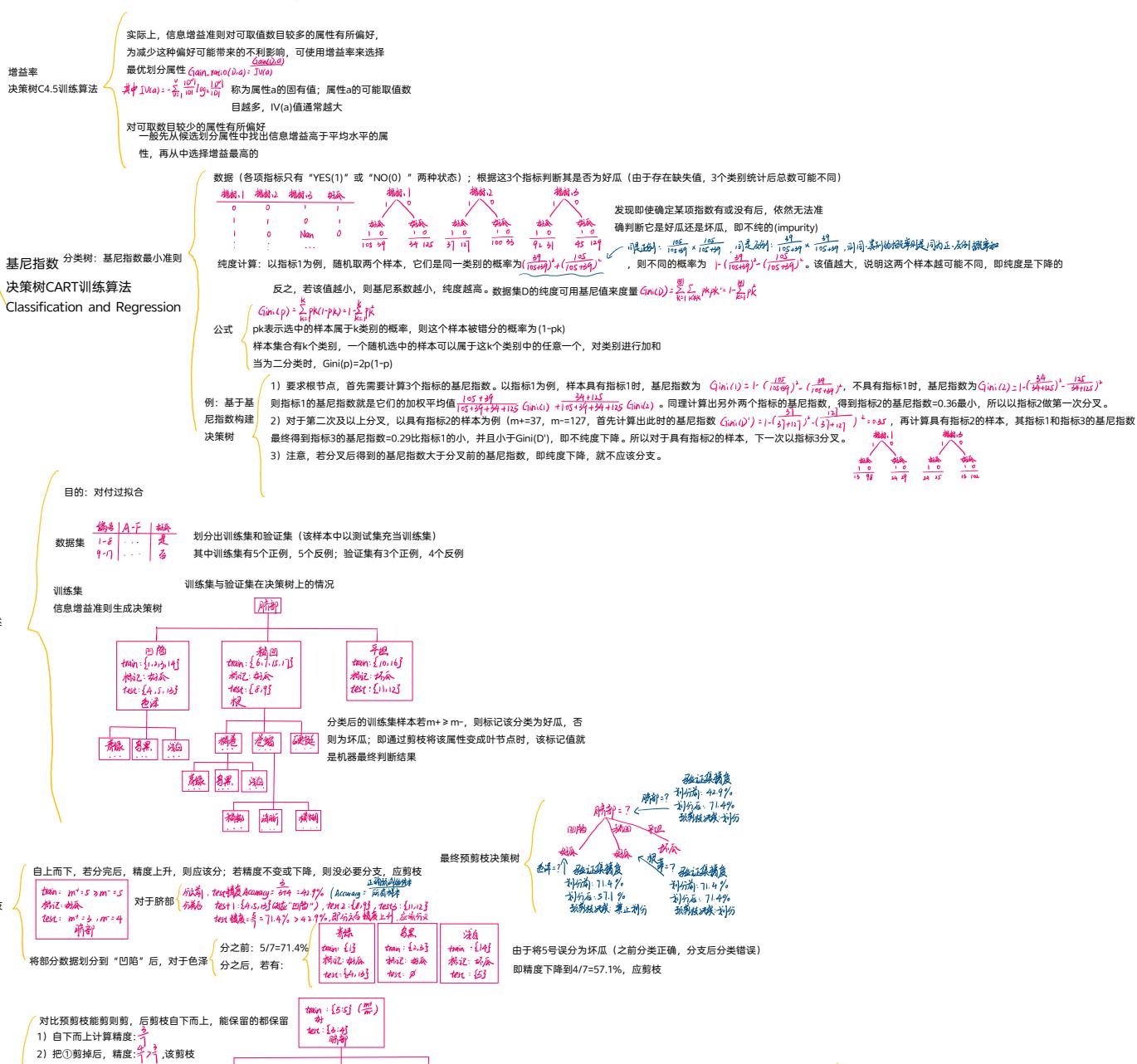


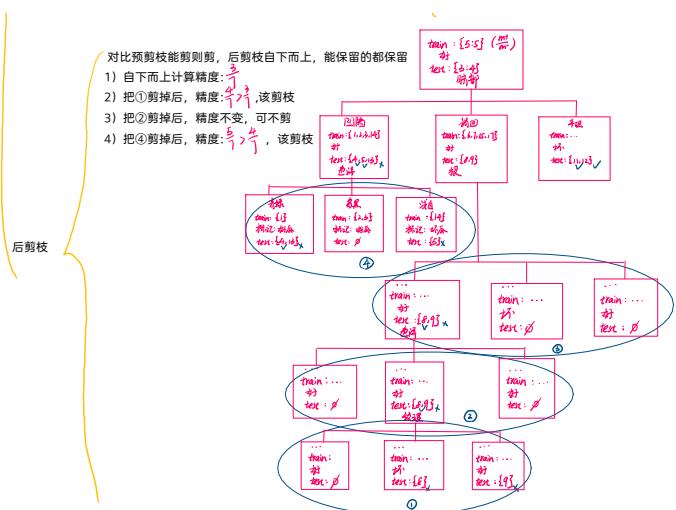
# 决策树



## 划分选择







连续与缺失值

连续值处理 二分法，使得信息增益最大的区间 均为连续值  
例：根据西瓜数据集对西瓜进行分类分析  
1) 把西瓜数据集从小到大排序并把数据样本排好  
2) 通过计算信息增益值，查找每一次分类的最佳区间。  
继续以剪枝处理的西瓜数据集D为例，即数据集为  

好	A	F	坏
1-6	...	...	是

  
其中 1-6是正例 0-7是反例

继续以剪枝处理的西瓜数据集口  
其中，1-8号正例，9-17号反例

缺失值处理

- 1) 以色泽属性为例, 缺失了1号, 5号, 13号在色泽上的数据, 即无缺失值的样例子集 $D_1$ 中只有14个样例, 计算分类前的信息熵 $Ent(D) = \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \approx 0.88$
- 2) 令 $\hat{D}_1, \hat{D}_2, \hat{D}_3$ 分别表示在色泽属性上取值为青绿、乌黑和浅白的样本子集, 有 $Ent(\hat{D}_1) = \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \approx 0.88$ ,  $Ent(\hat{D}_2) = \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \approx 0.88$ ,  $Ent(\hat{D}_3) = \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \approx 0.88$
- 3) 因此, 样本子集 $\hat{D}_1$ 在色泽的信息增益为 $Gain(D_1, \text{色泽}) = Ent(D) - Ent(\hat{D}_1) = \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 0.88$
- 4) 于是, 样本子集 $\hat{D}_1$ 上色泽属性的信息增益为 $Gain(\hat{D}_1, \text{色泽}) = \frac{1}{14} \cdot Gain(D_1, \text{色泽}) = \frac{1}{14} \cdot 0.88 \approx 0.06$ , 类似地可计算出其它属性在D上的信息增益

多变量决策树

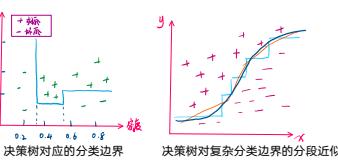
对于正常基于连续值生成的决策树，它们不同的属性最后往往呈线性的关系，因为每次的取值都是一次二分类

名錄標註範例

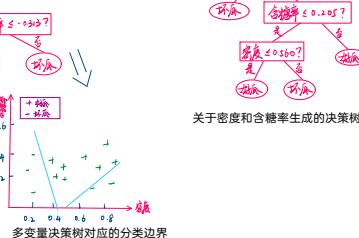
```

graph TD
    A["密度 ≤ 0.44?"] -- 否 --> B["含糖量 ≤ 0.245?"]
    A -- 是 --> C["含糖量 ≤ 0.245?"]
    B -- 否 --> D["密度 ≤ 0.65?"]
    B -- 是 --> E["密度 ≤ 0.65?"]
    C -- 否 --> F["含糖量 ≤ 0.245?"]
    C -- 是 --> G["含糖量 ≤ 0.245?"]
    D -- 否 --> H["密度 ≤ 0.65?"]
    D -- 是 --> I["密度 ≤ 0.65?"]
    E -- 否 --> J["密度 ≤ 0.65?"]
    E -- 是 --> K["密度 ≤ 0.65?"]
    F -- 否 --> L["密度 ≤ 0.65?"]
    F -- 是 --> M["密度 ≤ 0.65?"]
    G -- 否 --> N["密度 ≤ 0.65?"]
    G -- 是 --> O["密度 ≤ 0.65?"]
    H -- 否 --> P["密度 ≤ 0.65?"]
    H -- 是 --> Q["密度 ≤ 0.65?"]
    I -- 否 --> R["密度 ≤ 0.65?"]
    I -- 是 --> S["密度 ≤ 0.65?"]
    J -- 否 --> T["密度 ≤ 0.65?"]
    J -- 是 --> U["密度 ≤ 0.65?"]
    K -- 否 --> V["密度 ≤ 0.65?"]
    K -- 是 --> W["密度 ≤ 0.65?"]
    L -- 否 --> X["密度 ≤ 0.65?"]
    L -- 是 --> Y["密度 ≤ 0.65?"]
    M -- 否 --> Z["密度 ≤ 0.65?"]
    M -- 是 --> AA["密度 ≤ 0.65?"]
    N -- 否 --> BB["密度 ≤ 0.65?"]
    N -- 是 --> CC["密度 ≤ 0.65?"]
    O -- 否 --> DD["密度 ≤ 0.65?"]
    O -- 是 --> EE["密度 ≤ 0.65?"]
    P -- 否 --> FF["密度 ≤ 0.65?"]
    P -- 是 --> GG["密度 ≤ 0.65?"]
    Q -- 否 --> HH["密度 ≤ 0.65?"]
    Q -- 是 --> II["密度 ≤ 0.65?"]
    R -- 否 --> JJ["密度 ≤ 0.65?"]
    R -- 是 --> KK["密度 ≤ 0.65?"]
    S -- 否 --> LL["密度 ≤ 0.65?"]
    S -- 是 --> MM["密度 ≤ 0.65?"]
    T -- 否 --> NN["密度 ≤ 0.65?"]
    T -- 是 --> OO["密度 ≤ 0.65?"]
    U -- 否 --> PP["密度 ≤ 0.65?"]
    U -- 是 --> QQ["密度 ≤ 0.65?"]
    V -- 否 --> RR["密度 ≤ 0.65?"]
    V -- 是 --> SS["密度 ≤ 0.65?"]
    W -- 否 --> TT["密度 ≤ 0.65?"]
    W -- 是 --> YY["密度 ≤ 0.65?"]
    X -- 否 --> WW["密度 ≤ 0.65?"]
    X -- 是 --> ZZ["密度 ≤ 0.65?"]
    Y -- 否 --> XX["密度 ≤ 0.65?"]
    Y -- 是 --> YY["密度 ≤ 0.65?"]
    AA -- 否 --> BB["密度 ≤ 0.65?"]
    AA -- 是 --> CC["密度 ≤ 0.65?"]
    BB -- 否 --> DD["密度 ≤ 0.65?"]
    BB -- 是 --> EE["密度 ≤ 0.65?"]
    CC -- 否 --> FF["密度 ≤ 0.65?"]
    CC -- 是 --> GG["密度 ≤ 0.65?"]
    DD -- 否 --> HH["密度 ≤ 0.65?"]
    DD -- 是 --> II["密度 ≤ 0.65?"]
    EE -- 否 --> JJ["密度 ≤ 0.65?"]
    EE -- 是 --> KK["密度 ≤ 0.65?"]
    FF -- 否 --> LL["密度 ≤ 0.65?"]
    FF -- 是 --> MM["密度 ≤ 0.65?"]
    GG -- 否 --> NN["密度 ≤ 0.65?"]
    GG -- 是 --> OO["密度 ≤ 0.65?"]
    HH -- 否 --> PP["密度 ≤ 0.65?"]
    HH -- 是 --> QQ["密度 ≤ 0.65?"]
    II -- 否 --> RR["密度 ≤ 0.65?"]
    II -- 是 --> SS["密度 ≤ 0.65?"]
    JJ -- 否 --> TT["密度 ≤ 0.65?"]
    JJ -- 是 --> YY["密度 ≤ 0.65?"]
    KK -- 否 --> WW["密度 ≤ 0.65?"]
    KK -- 是 --> ZZ["密度 ≤ 0.65?"]
    LL -- 否 --> XX["密度 ≤ 0.65?"]
    LL -- 是 --> YY["密度 ≤ 0.65?"]
    MM -- 否 --> YY["密度 ≤ 0.65?"]
    MM -- 是 --> ZZ["密度 ≤ 0.65?"]
    NN -- 否 --> YY["密度 ≤ 0.65?"]
    NN -- 是 --> YY["密度 ≤ 0.65?"]
    OO -- 否 --> YY["密度 ≤ 0.65?"]
    OO -- 是 --> YY["密度 ≤ 0.65?"]
    PP -- 否 --> YY["密度 ≤ 0.65?"]
    PP -- 是 --> YY["密度 ≤ 0.65?"]
    QQ -- 否 --> YY["密度 ≤ 0.65?"]
    QQ -- 是 --> YY["密度 ≤ 0.65?"]
    RR -- 否 --> YY["密度 ≤ 0.65?"]
    RR -- 是 --> YY["密度 ≤ 0.65?"]
    SS -- 否 --> YY["密度 ≤ 0.65?"]
    SS -- 是 --> YY["密度 ≤ 0.65?"]
    TT -- 否 --> YY["密度 ≤ 0.65?"]
    TT -- 是 --> YY["密度 ≤ 0.65?"]
    YY -- 否 --> YY["密度 ≤ 0.65?"]
    YY -- 是 --> YY["密度 ≤ 0.65?"]
    ZZ -- 否 --> YY["密度 ≤ 0.65?"]
    ZZ -- 是 --> YY["密度 ≤ 0.65?"]

```



多变量决策树的分界点不某个属性，而是对属性的会进行一种测试



多变量决策树的分界点不再是对于某个属性，而是对属性的线性组合进行一种测试

多变量决策树对应的分类边界

