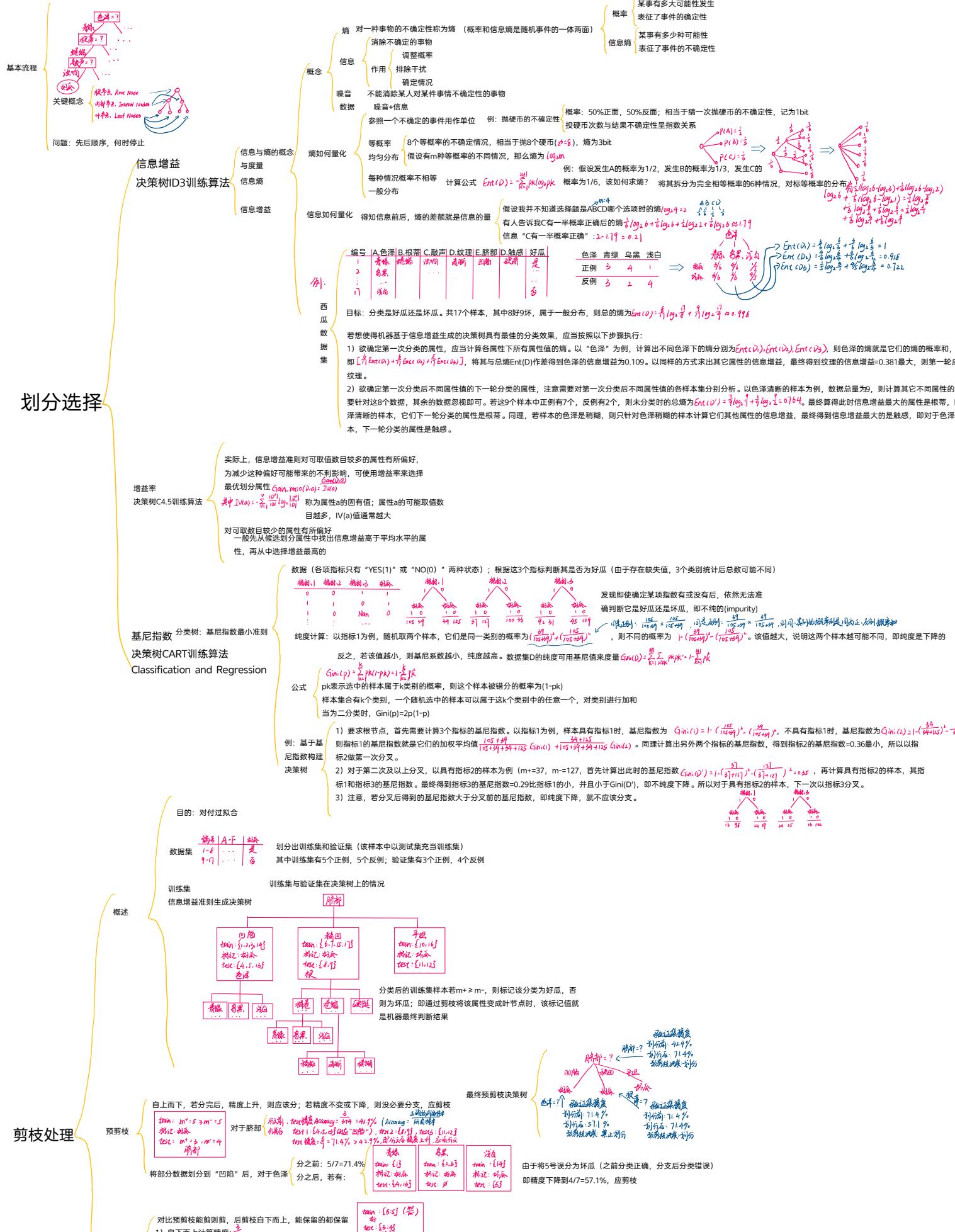
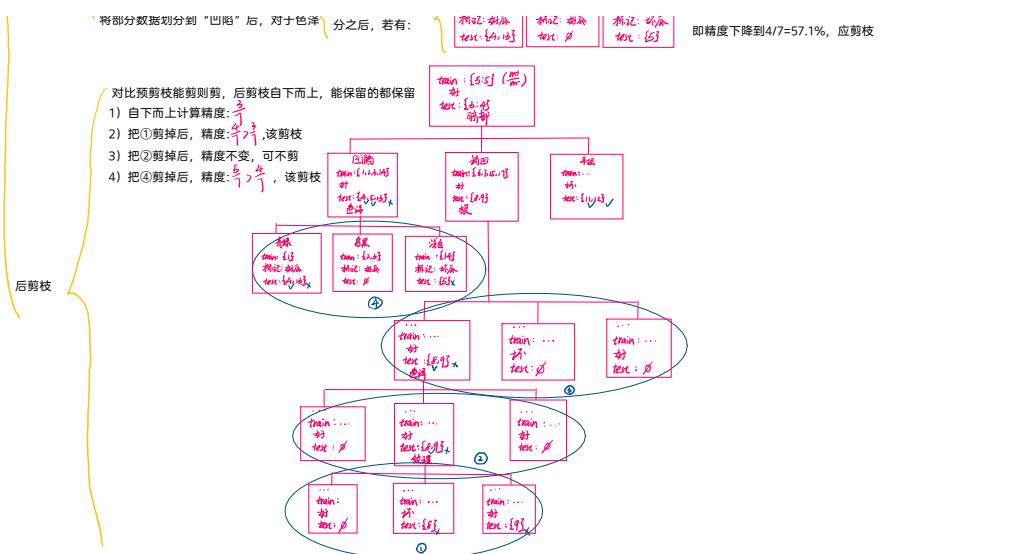


决策树





连续与缺失值

连续值处理 二分法，使得信息增益最大的阈值 均为连续值

例：数据集 D 对西瓜数据集分类

- 1) 按照次序从小到大的顺序把数据集样本排序
- 2) 通过计算信息增益，查找每一次分类的最佳阈值。

继续以剪枝处理的西瓜数据集 D 为例，即数据集为 $D = \{(1, \text{正}), (2, \text{正}), (3, \text{正}), (4, \text{正}), (5, \text{正}), (6, \text{正}), (7, \text{正}), (8, \text{正}), (9, \text{反}), (10, \text{反}), (11, \text{反}), (12, \text{反}), (13, \text{反})\}$ ，其中，1-8 是正例，9-13 是反例。

- 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' 中只有 14 个样例，计算分类前的信息熵 $Ent(D) = \frac{1}{14} \log_2 \frac{1}{14} + \frac{1}{14} \log_2 \frac{1}{14} \approx 0.94$
- 2) 令 $\text{v}_1, \text{v}_2, \text{v}_3$ 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 $Ent(\text{v}_1) = \frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \approx 1.0$, $Ent(\text{v}_2) = \frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918$, $Ent(\text{v}_3) = \frac{1}{3} \log_2 \frac{1}{3} = 0$
- 3) 因此，样本子集 D' 上色泽的信息增益为 $Gain(D', \text{色泽}) = Ent(D) - \frac{1}{14} [Ent(\text{v}_1) + Ent(\text{v}_2) + Ent(\text{v}_3)] = 0.09$
- 4) 于是，样本子集 D 上色泽属性的信息增益为 $Gain(D, \text{色泽}) = Gain(D', \text{色泽}) \cdot \frac{14}{17} = 0.053$ ，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

