Lecture 6: Decision Trees

Madhavan Mukund

https://www.cmi.ac.in/~madhavan

Data Mining and Machine Learning August-December 2020

Training data

- Each item is characterized by attributes (a_i, a₂,..., a_k)
- Each item is assigned a class or category c

Goal

Predict *c* for a new item with attributes $(a'_1, a'_2, \ldots, a'_k)$

ID	Age	Has_job	Own_house	Credit_rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No

Decision trees

- Play "20 Questions" with the training data
- Query an attribute
 - Partition the training data based on the answer
- Repeat until you reach a partition with a uniform category
- Queries are adaptive
 - Different along each path, depends on history



ID	Age	Has_job	Own_house	Credit_rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
16	old	falaa	falsa	fair	Ne

A : current set of attributes

Pick $a \in A$, create children corresponding to resulting partition with attributes $A \setminus \{a\}$

Stopping criterion:

- Current node has uniform class label
- A is empty no more attributes to query

If a leaf node is not uniform, use majority class as prediction



- Non-uniform leaf node identical combination of attributes, but different classes
- Attributes do not capture all criteria used for classification

Decision trees

- Tree is not unique
- Which tree is better?
- Prefer small trees
 - Explainability
 - Generalize better (see later)

Unfortunately

- Finding smallest tree is NP-complete — for any definition of "smallest"
- Instead, greedy heuristic





Greedy heuristic

- Goal: partition with uniform category — pure leaf
- Impure node best prediction is majority value
- Minority ratio is impurity
- Heuristic: reduce impurity as much as possible
- For each attribute, compute weighted average impurity of children
- Choose the minimum
- Will see better heuristics



