

Data Mining Tutorial

Mark A. Austin

University of Maryland

austin@umd.edu

ENCE 688P, Fall Semester 2021

October 16, 2021

Mathematical Framework

Information Gain

The amount of information that is gained by knowing the value of an attribute. It equals the entropy of a distribution before a split minus the entropy of a distribution after a split.

$$IG(Y, X) = H(Y) - H(Y|X). \quad (17)$$

Here:

- Information gain $IG(X, Y)$ is the reduction of uncertainty about Y given an additional piece of information X about Y .
- $H(Y)$ is the entropy of Y (before split).
- $H(Y|X)$ is the conditional entropy of Y given the value of attribute X (after split).

Example 1 (Buy Computer)

Initial Dataset. Will customer buy a computer?

ID	Age Group	Income	Student	Credit Rating	Buys Computer
1	young	high	no	fair	no
2	young	high	no	excellent	no
3	middle	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle	low	yes	excellent	yes
8	young	medium	no	fair	no
9	young	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	young	medium	yes	excellent	yes
12	middle	medium	no	excellent	yes
13	middle	high	yes	fair	yes
14	senior	medium	no	excellent	no

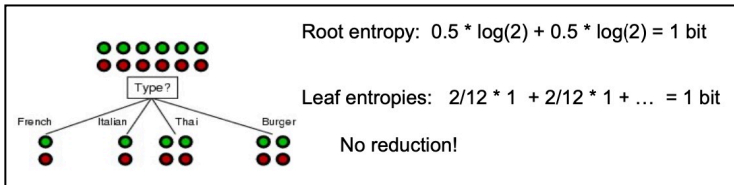
Example 2 (Customer Wait for Table at Restaurant?)

Dataset Attributes

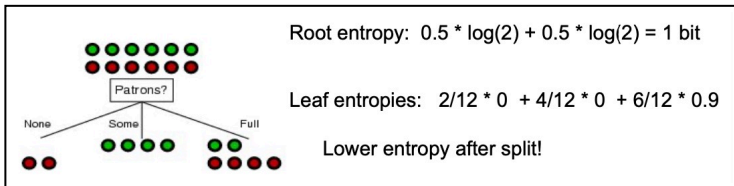
- **Alternate:** Is there a suitable alternate restaurant nearby?
- **Bar:** Does restaurant have comfortable bar area to wait in?
- **Fri/Sat:** True on Fridays and Saturdays.
- **Hungry:** True when customer is hungry.
- **Patrons:** How many people are in the restaurant? (none, some, and full).
- **Price:** The restaurant price range (\$, \$\$ and \$\$\$).
- **Raining:** Is it raining outside?
- **Reservation:** Did customer make a reservation?
- **Type:** Type of restaurant (French, Italian, Thai, or Burger).
- **WaitEstimate:** Wait time estimated by host (0-10 mins, 10-30, 30-60, or > 60).

Example 2 (Customer Wait for Table at Restaurant?)

Split on Restaurant Type Attribute

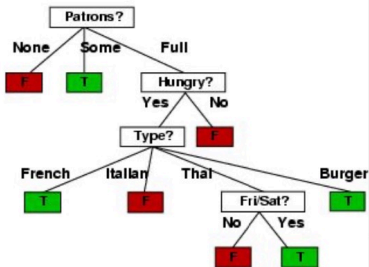
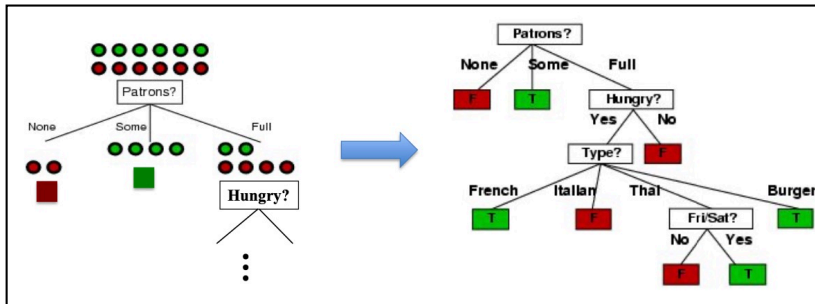


Split on Patrons Attribute



Example 2 (Customer Wait for Table at Restaurant?)

Decision Tree Synthesis



Classification with Decision Trees (Summary)

Advantages

- Decision trees are simple to understand and interpret.
- Requires only a small number of observations.
- Best and expected values can be determined for different scenarios.

Disadvantages

- Difficulties in handling data with missing values.
- Information gain criterion is biased in favor of attributes with more levels.
- Calculations become complex if values are uncertain or outcomes are linked.

References

- Jaynes E.T., Information Theory and Statistical Mechanics. II, Phys. Rev. 108, 171, October 1957.
- Kapur J.N., Maximum-Entropy Models in Science and Engineering, John Wiley and Sons, 1989.
- Mitchell T.M., Machine Learning and Data Mining, Communications of the ACM, Vol. 42., No. 11, November 1999.
- Russell S., and Norvig P., Artificial Intelligence: A Modern Approach (Third Edition), Prentice-Hall, 2010.
- Shanon C.E., and Weaver W., The Mathematical Theory of Communication, University of Illinois, Urbana, Chicago, 1949.
- Witten I.H., Frank E., Hall M.A., and Pal C.J., Data Mining: Practical Machine Learning Tools and Techniques, Morgan Kaufmann, 2017.