

Test #2

CS5402 — Intro To Data Mining

Illya Starikov

Due Date: July 29th, 2018

1 C4.5

$$Split = 120$$

$$Entropy_{LT} = -\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right)$$

$$Entropy_{Gt} = -\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right)$$

$$Info\ Gain = X - \left(\frac{4}{9} \times Entropy_{LT} + \frac{5}{9} \times Entropy_{GT} \right)$$

$$Split = 140$$

$$Entropy_{LT} = -\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right)$$

$$Entropy_{Gt} = -\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right)$$

$$Info\ Gain = X - \left(\frac{6}{9} \times Entropy_{LT} + \frac{3}{9} \times Entropy_{GT} \right)$$

We choose the split that we calculated previously based on it's associated **highest information gain**. Values before it should be assigned *Less Than Split* and values after it should be assigned *Greater Than Split*.

2 Grouping Or Splitting

$$\{\{SciFiction, Mystery\}, \{NonFiction\}\}$$

$$SciMystery = -\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) - \frac{1}{6} \log_2 \left(\frac{1}{6} \right)$$

$$NonFiction = -0 - \frac{1}{6} \log_2 \left(\frac{1}{6} \right) - \frac{1}{6} \log_2 \left(\frac{1}{6} \right)$$

$$Info\ Gain = 1.5575 - \left(\frac{6}{8} \times SciMystery + \frac{2}{8} \times NonFiction \right)$$

$$\{\{SciFiction, NonFiction\}, \{Mystery\}\}$$

$$SciNonFiction = -\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right)$$

$$NonFiction = -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right)$$

$$Info\ Gain = 1.5575 - \left(\frac{4}{8} \times SciNonFiction + \frac{4}{8} \times Mystery \right)$$

$$\{\{Mystery, NonFiction\}, \{NonFiction\}\}$$

$$\text{SciNonFiction} = -1/6 \log_2 (1/6) - 3/6 \log_2 (3/6) - 2/6 \log_2 (2/6)$$

$$\text{NonFiction} = -2/2 \log_2 (2/2) - 0 - 0$$

$$\text{Info Gain} = 1.5575 - (6/8 \times \text{SciNonFiction} + 2/8 \times \text{NonFiction})$$

To determine if any attributes are to be grouped, take the **highest information gain**, and compare it to the Entropy Before Split for music preference. If information gain is higher, then the grouping was better; if not, the grouping was worse.

3 Support Vectors

The equations would reduce down to:

$$10 \alpha_1 + 10 \alpha_2 = -1$$

$$10 \alpha_1 + 11 \alpha_2 = 1$$

For which we get the following solutions:

$$\alpha_1 = -\frac{21}{10} \quad \alpha_2 = 2$$

For the discriminating 2D hyperplane, we get as follows:

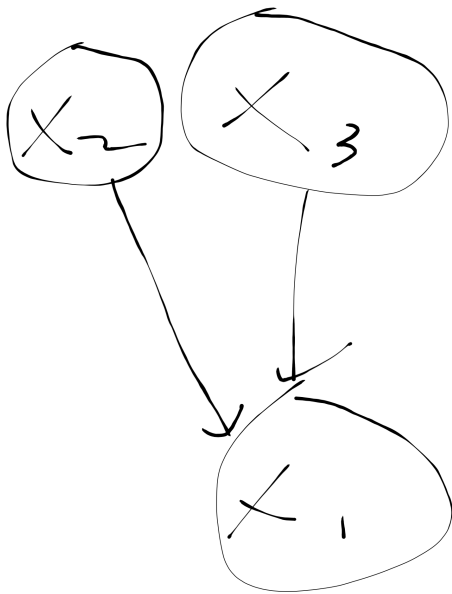
$$\begin{aligned} w &= -\frac{21}{10} \times s_1 + 2 \times s_2 \\ &= \left\langle -\frac{3}{10}, 2, -\frac{41}{10} \right\rangle \end{aligned}$$

Meaning our equation would be as follows:

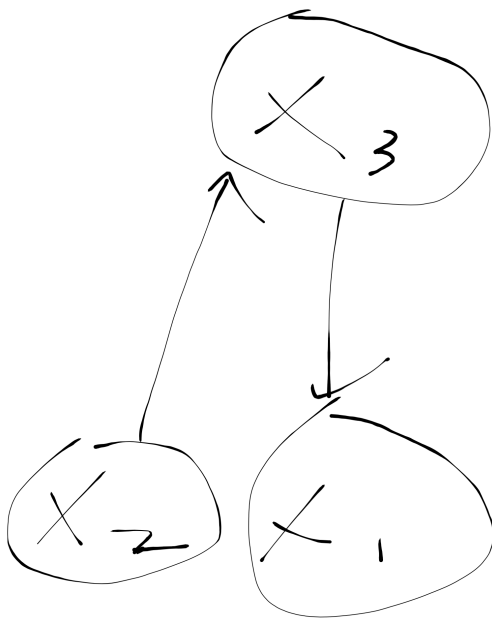
$$-3/10 x - 2 y - 41/10 = 0$$

4 Bayesian Network

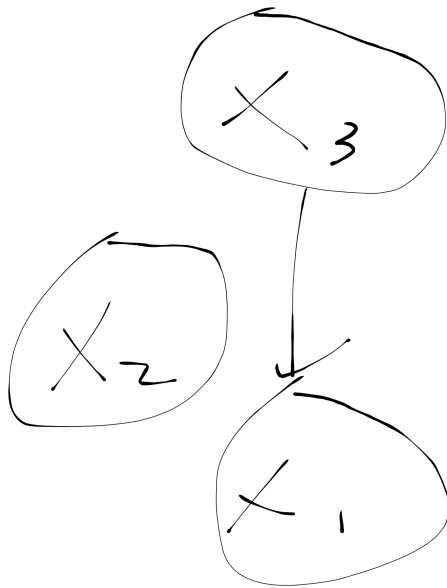
The algorithm will enumerate the possible options (in the original ordering) for a parent of X2. It does so by calculating the probability of the particular attribute being a parent of X2, and the one with the greatest probability is the node that becomes the parents.



1.



2.



3.

5 DBScan

	density	designation
A1	2	Border
A2	2	Border
A3	3	Core
A4	1	Noise
A5	2	Border
A6	2	Border
A7	2	Noise
A8	1	Noise
A9	2	Noise
A10	3	Core

The clusters would be as follows: $\{A1, A2, A3\}$ and $\{A5, A6, A10\}$. Points with have are dense (density > number of points) are designated as core points. Points that are around a core points (within an ϵ of a core point) are designated as border point. Points that fall in neither of these categories are designated as noise.

6 Linkage

The single linkage would be as follows:

$$|(1, 4) - (3, 5)| = 3$$

The complete linkage would be as follows:

$$|(1, 2) - (4, 5)| = 6$$

The centroid linkage would be as follows:

$$|(1, 11/4) - (14/4, 5)| = 19/4$$

The average linkage would be as follows:

$$\text{Average of all distances between points} = 5$$

7 Ensemble Classifier

Sum	-2	4	-6	2
Class	-1	1	-1	1

8 Bayes Network

For **worker** = **T**, the likelihood would be as follows:

$$0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.01$$

For **worker** = **F**, the likelihood would be as follows:

$$0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.99$$

The probability for **worker** = **T** would be:

$$\frac{0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.01}{0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.01 + 0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.99}$$

And the probability for **worker** = **F** would be:

$$\frac{0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.99}{0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.01 + 0.8 \times 0.02 \times 0.1 \times 0.1 \times 0.99}$$

Multiple Choice

9. d. repeatedly sampling from the original dataset according to a uniform probability distribution
10. c. increased, decreased

- 11. a. bias, variance
- 12. b. false
- 13. c. 180 instances reached this point in the decision tree, but 22 of those were not classified as `tested_negative` in the training dataset
- 14. d. $\text{age} > 34$
- 15. b. compute the predicted error rate of the rule with one condition (C_1 , C_2 , C_3) deleted and no conditions deleted, and, from those, use the version of the rule with the lowest rate
- 16. c. each instance in the dataset will have a probability of being in a particular cluster
- 17. b. supervised