

CSCI 567: Theory Assignment 1

Anurima Anil Padwal
USC ID: 4348819703

September 2019

Solutions

1. For $\mathbf{x}_i, \mathbf{x}_o, \mathbf{x}_j$ normalized to the unit norm, $\|\mathbf{x}_i\|_2 = \|\mathbf{x}_o\|_2 = \|\mathbf{x}_j\|_2 = 1$. The cosine distance between \mathbf{x}_i and \mathbf{x}_j is

$$\begin{aligned} C(\mathbf{x}_i, \mathbf{x}_j) &= 1 - \frac{\mathbf{x}_i^T \cdot \mathbf{x}_j}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2} \\ &= 1 - \mathbf{x}_i^T \cdot \mathbf{x}_j \end{aligned} \quad (1)$$

Similarly, the cosine distance between \mathbf{x}_i and \mathbf{x}_o

$$\begin{aligned} C(\mathbf{x}_i, \mathbf{x}_o) &= 1 - \frac{\mathbf{x}_i^T \cdot \mathbf{x}_o}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_o\|_2} \\ &= 1 - \mathbf{x}_i^T \cdot \mathbf{x}_o \end{aligned} \quad (2)$$

Euclidian distance between \mathbf{x}_i and \mathbf{x}_j is given by:

$$\begin{aligned} E(\mathbf{x}_i, \mathbf{x}_j) &= \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \\ &= (\mathbf{x}_i - \mathbf{x}_j)^T \cdot (\mathbf{x}_i - \mathbf{x}_j) \\ &= \mathbf{x}_i^T \mathbf{x}_j - 2 \cdot \mathbf{x}_i^T \mathbf{x}_j^T + \mathbf{x}_j^T \mathbf{x}_j \\ &= 1 + 1 - 2 \cdot \mathbf{x}_i^T \mathbf{x}_j \\ &= 2 - 2 \cdot \mathbf{x}_i^T \mathbf{x}_j \\ &= 2(1 - \mathbf{x}_i^T \mathbf{x}_j) \end{aligned} \quad (3)$$

Similarly, Euclidian distance between \mathbf{x}_i and \mathbf{x}_o is given by:

$$\begin{aligned} E(\mathbf{x}_i, \mathbf{x}_o) &= \|\mathbf{x}_i - \mathbf{x}_o\|_2^2 \\ &= (\mathbf{x}_i - \mathbf{x}_o)^T \cdot (\mathbf{x}_i - \mathbf{x}_o) \\ &= \mathbf{x}_i^T \mathbf{x}_o - 2 \cdot \mathbf{x}_i^T \mathbf{x}_o^T + \mathbf{x}_j^T \mathbf{x}_o \\ &= 1 + 1 - 2 \cdot \mathbf{x}_i^T \mathbf{x}_o \\ &= 2 - 2 \cdot \mathbf{x}_i^T \mathbf{x}_o \\ &= 2(1 - \mathbf{x}_i^T \mathbf{x}_o) \end{aligned} \quad (4)$$

Given,

$$C(\mathbf{x}_i \cdot \mathbf{x}_j) \leq C(\mathbf{x}_i \cdot \mathbf{x}_o)$$

Multiplying on both sides by 2,

$$2C(\mathbf{x}_i \cdot \mathbf{x}_j) \leq 2C(\mathbf{x}_i \cdot \mathbf{x}_o)$$

From (1) and (2),

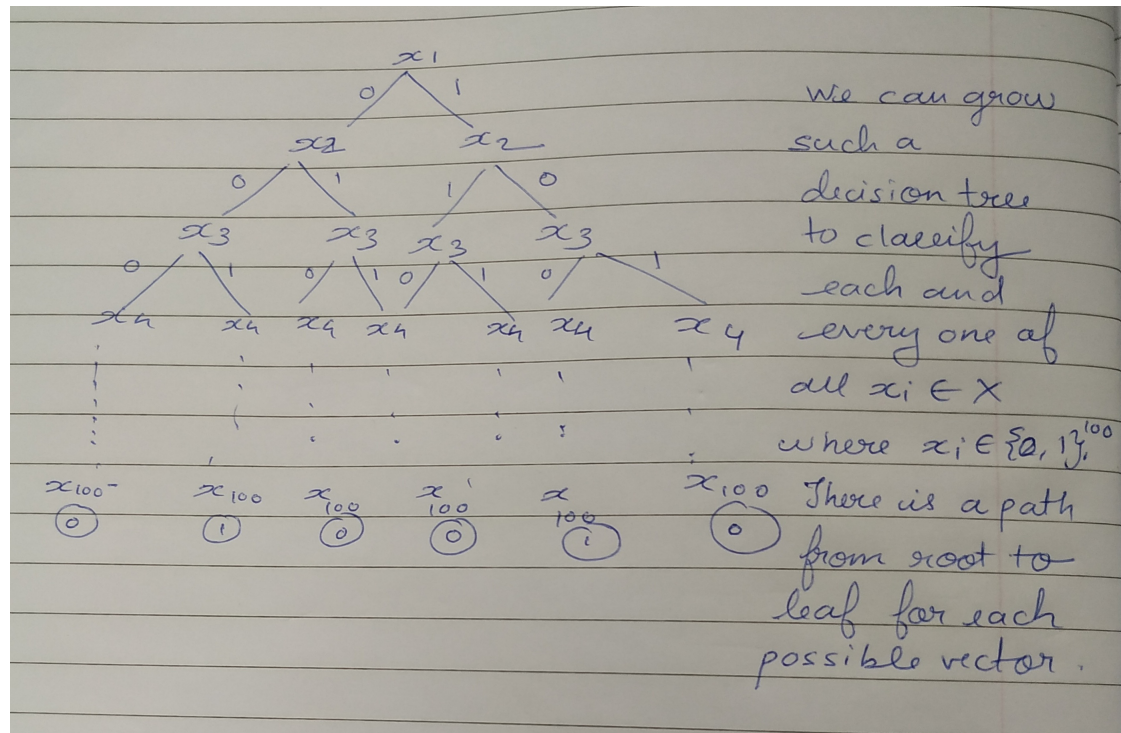
$$2(1 - \mathbf{x}_i^T \cdot \mathbf{x}_j) \leq 2(1 - \mathbf{x}_i^T \cdot \mathbf{x}_o)$$

From (3) and (4),

$$E(\mathbf{x}_i, \mathbf{x}_j) \leq E(\mathbf{x}_i, \mathbf{x}_o)$$

Hence, proved.

2.1 Yes, we can have a decision tree to classify a dataset of 100 dimensional vectors, such that each vector is a binary vector $\mathbf{x} \in \{0, 1\}^{100}$, and the tree has zero classification error. First we generate all possible combinations for each such \mathbf{x}_i . There are 2^{100} such possibilities since the dataset has binary vectors of 100 dimensions. Here, we split on any feature in the vector for the root node, and continue building the tree top down, such that we have a path from root to leaf node for each \mathbf{x}_i .



2.2 Yes. First we generate all possible combinations for each such \mathbf{x}_i . There are 2^{100} such possibilities since the dataset has binary vectors of 100 dimensions. We can classify this entire set using our decision tree, and use the resulting set of vectors and labels as our 1-NN classifier. This will give the same result as the decision tree since each point is its own neighbor.

3. Yes, the given decision tree can be implemented as a 1 NN classifier. The values are as follows:

(x_1, x_2)	Label
(A + 1, B + 1)	1
(A + 1, B - 1)	0
(A - 1, B + 1)	0
(A - 1, B - 1)	1

4.1 Based on the decision tree in Figure 3, the number of mis-classifications is 0. Hence the test error = $\frac{\text{number of misclassifications}}{\text{number of total test samples}} = \frac{0}{2} = 0$

4.2 Based on the decision tree in Figure 4, the number of misclassifications on the test data is 1. The point $(x_1 = 0.2, x_2 = 0.8)$ is misclassified. Hence the test error = $\frac{\text{number of misclassifications}}{\text{number of total test samples}} = \frac{1}{2} = 0.5$

4.3 Yes, the decision tree in Figure 4 is a linear classifier in terms of (x_1, x_2) . The classifier may be represented as $y = 1$ when $z > 0.5$ and $y = 0$ for all other cases. Here z represents the weighted combination $z = w_1 \cdot x_1 + w_2 \cdot x_2$ where $w_2 = 0$ and $w_1 \geq 1$. No, we cannot classify the data and get zero classification error by drawing a depth 1 decision tree.

4.4 No.

Justification:

Putting the data in the rectangle:

$$ax_1 + bx_2 \geq c$$

$$ax_1 + bx_2 < c$$

we get,

$$a \geq 1 \quad (1)$$

$$b \geq 1 \quad (2)$$

$$0 < 1 \quad (3)$$

$$a + b < 1 \quad (4)$$

Equations (1), (2), and (4) are contradictory. Hence, the data in Table 1 is not linearly separable. Hence, we cannot classify the data using a depth 1 decision tree to get a zero classification error.

5. The structure of T_1 is:

Left child:

class A: 150 samples

class B: 50 samples

Label: class A

Right child:

class A: 50 samples

class B: 150 samples

Label: class B

The structure of T_2 is:

Left child:

class A: 0 samples

class B: 100 samples

Label: class B

Right child:

class A: 200 samples

class B: 100 samples

Label: class A

5.1 For tree T_1 ,

Left child:

50 samples are misclassified as class B.

Classification error = $\frac{50}{200} = 0.25$

Right child:

50 samples are misclassified as class A.

Classification error = $\frac{50}{200} = 0.25$

For tree T_2 ,

Left child:

0 samples are misclassified

Classification error = $\frac{0}{100} = 0$

Right child:

100 samples are misclassified as class B.

Classification error = $\frac{100}{300} = 0.33$

Entropy is given by,

$$H(P) = - \sum_{k=1}^C P(Y = k) \log P(Y = k)$$

For tree T_1 ,

Left child:

$$\begin{aligned}
 \text{Entropy} &= -\frac{150}{150+50} \log_e \frac{150}{150+50} - \frac{50}{150+50} \log_e \frac{50}{150+50} \\
 &= -\frac{150}{200} \log_e \frac{150}{200} - \frac{50}{200} \log_e \frac{50}{200} \\
 &= 0.56
 \end{aligned}$$

Right child:

$$\begin{aligned}
 \text{Entropy} &= -\frac{150}{150+50} \log_e \frac{150}{150+50} - \frac{50}{150+50} \log_e \frac{50}{150+50} \\
 &= -\frac{150}{200} \log_e \frac{150}{200} - \frac{50}{200} \log_e \frac{50}{200} \\
 &= 0.56
 \end{aligned}$$

For tree T_2 ,

Left child:

$$\begin{aligned}
 \text{Entropy} &= -\frac{0}{100+0} \log_e \frac{0}{100+0} - \frac{100}{100+0} \log_e \frac{100}{100+0} \\
 &= 0
 \end{aligned}$$

Right child:

$$\begin{aligned}
 \text{Entropy} &= -\frac{200}{200+100} \log_e \frac{200}{200+100} - \frac{100}{200+100} \log_e \frac{100}{200+100} \\
 &= -\frac{200}{300} \log_e \frac{200}{300} - \frac{100}{300} \log_e \frac{100}{300} \\
 &= 0.64
 \end{aligned}$$

Gini Impurity

For tree T_1 ,

Left child,

$$\begin{aligned}
 \text{Gini Impurity} &= 1 - \sum p_i^2 \\
 &= 1 - \left(\frac{150}{200}\right)^2 - \left(\frac{50}{200}\right)^2 \\
 &= 0.38
 \end{aligned}$$

Right child,

$$\begin{aligned}
 \text{Gini Impurity} &= 1 - \sum p_i^2 \\
 &= 1 - \left(\frac{150}{200}\right)^2 - \left(\frac{50}{200}\right)^2 \\
 &= 0.38
 \end{aligned}$$

For tree T_2 , Left child,

$$\begin{aligned}\text{Gini Impurity} &= 1 - \sum p_i^2 \\ &= 1 - \left(\frac{0}{100}\right)^2 - \left(\frac{100}{100}\right)^2 \\ &= 0\end{aligned}$$

Right child,

$$\begin{aligned}\text{Gini Impurity} &= 1 - \sum p_i^2 \\ &= 1 - \left(\frac{200}{300}\right)^2 - \left(\frac{100}{300}\right)^2 \\ &= 0.44\end{aligned}$$

5.2 T_1 misclassifies 100 samples in all and T_2 also misclassifies 100 samples in all. Hence we cannot say one is better than the other in terms of classification error as it is the same of both, i.e. 0.25.

The conditional entropy of T_1 is:

$$\begin{aligned}H_1 * P_1 + H_2 * P_2 \\ = 0.56 * \frac{200}{400} + 0.56 * \frac{200}{400} = 0.56\end{aligned}$$

The conditional entropy of T_2 is:

$$\begin{aligned}H_1 * P_1 + H_2 * P_2 \\ = 0 * \frac{100}{400} + 0.63 * \frac{300}{400} = 0.48\end{aligned}$$

The conditional entropy of T_2 is less than that of T_1 . Hence, the split at T_2 is more pure than at T_1 .

Now, weighted Gini for $T_1 = \frac{200}{400} * 0.37 + \frac{200}{400} * 0.37 = 0.37$

Weighted Gini for $T_2 = 0 * \frac{100}{400} + 0.44 * \frac{300}{400} = 0.33$

Weighted Gini of T_2 is less than that of T_1 . Hence T_2 is more pure and has a better quality of split than T_1 .

$$6.1. P(\text{PlayTennis} = \text{Yes}) = \frac{4}{6} = \frac{2}{3} \quad P(\text{PlayTennis} = \text{No}) = \frac{2}{6} = \frac{1}{3}$$

$$6.2 \quad P(\text{Weather} = \text{Sunny} \mid \text{PlayTennis} = \text{Yes}) = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{Emotion} = \text{Normal} \mid \text{PlayTennis} = \text{Yes}) = \frac{1}{4}$$

$$P(\text{HomeWork} = \text{Much} \mid \text{PlayTennis} = \text{Yes}) = \frac{1}{4}$$

6.3

$$\begin{aligned}
& P(\text{PlayTennis} = \text{Yes} \mid \mathbf{x}) \\
&= \frac{P(\mathbf{x} \mid \text{PlayTennis} = \text{Yes}) * P(\text{PlayTennis} = \text{Yes})}{P(x)} \\
&= \frac{P(\frac{\text{Weather} = \text{Sunny}}{\text{Yes}}) * P(\frac{\text{Emotion} = \text{Normal}}{\text{Yes}}) * P(\frac{\text{Homework} = \text{Much}}{\text{Yes}})}{P(x)} \\
&= \frac{\frac{1}{2} * \frac{1}{4} * \frac{1}{4} * \frac{2}{3}}{P(x)} \\
&= \frac{1/48}{P(x)} \\
& P(\text{PlayTennis} = \text{No} \mid \mathbf{x}) \\
&= \frac{P(\mathbf{x} \mid \text{PlayTennis} = \text{No}) * P(\text{PlayTennis} = \text{No})}{P(x)} \\
&= \frac{P(\frac{\text{Weather} = \text{Sunny}}{\text{No}}) * P(\frac{\text{Emotion} = \text{Normal}}{\text{No}}) * P(\frac{\text{Homework} = \text{Much}}{\text{No}})}{P(x)} \\
&= \frac{\frac{1}{2} * \frac{1}{2} * \frac{1}{2} * \frac{1}{3}}{P(x)} \\
&= \frac{1/24}{P(x)}
\end{aligned}$$

The numerator of $P(\text{PlayTennis} = \text{No} \mid \mathbf{x})$ is greater than that of $P(\text{PlayTennis} = \text{Yes} \mid \mathbf{x})$. Hence $P(\text{PlayTennis} = \text{No} \mid \mathbf{x})$ has a larger value.