

11/17/24

CS483 - Big Data Mining HW4

Gokshulnath ①
Thirumaran
67-5086474

1) a) The impurity is

$$\begin{aligned}
 I(D) &= 100 \times (1 - (0.4)^2 - (0.6)^2) \\
 &= 100 \times (1 - 0.16 - 0.36) \\
 &= 100 \times 0.48 \\
 &= 48
 \end{aligned}$$

1.) So, 20 out of 50 wine drinkers like beer

20 out of 50 non-wine drinkers like beer.

Impurity of the "like wine" is $I(D_L)$

$$\begin{aligned}
 I(D_L) &= 50 \times (1 - (0.4)^2 - (0.6)^2) \\
 &= 50 \times (1 - 0.16 - 0.36) \\
 &= 50 \times 0.48
 \end{aligned}$$

$$I(D_L) = 24$$

Impurity of the "doesn't like wine" is $I(D_R)$

$$\begin{aligned}
 I(D_R) &= 50 \times (1 - (0.4)^2 - (0.6)^2) \\
 &= 50 \times 0.48
 \end{aligned}$$

$$I(D_R) = 24$$

(2)

Thus,

$$\begin{aligned}
 G_1 &= I(0) - (I(D_L) + I(D_R)) \\
 &= 48 - (24 + 24) \\
 &= 48 - 48 = 0.
 \end{aligned}$$

So no reduction in impurity.

2. So, 20 out of 30 runners like beer
20 out of 70 non-runners like beer

The impurity on "like running" side is $I(D_L)$

$$\begin{aligned}
 I(D_L) &= 30 \times (1 - (0.66)^2 - (0.33)^2) \\
 &= 30 \times (1 - 0.4356 - 0.1089) \\
 &= 30 \times 0.4555
 \end{aligned}$$

$$I(D_L) = 13.665$$

The impurity on "doesn't like running" is $I(D_R)$

$$\begin{aligned}
 I(D_R) &= 70 \times (1 - (0.2857)^2 - (0.7143)^2) \\
 &= 70 \times (1 - (0.08162) - (0.51022)) \\
 &= 70 \times (0.40815)
 \end{aligned}$$

$$I(D_R) = 28.5714$$

Thus,

$$\begin{aligned} G_1 &= I(D_1) - (I(D_L) + I(D_R)) \\ &= 48 - (13.665 + (28.5714)) \\ &= 48 - 42.2364 \\ &= 5.7636 \end{aligned}$$

3) So, 50 out of 80 pizza lovers like beer
10 out of 20 non-pizza lovers like beer

The impurity on "like pizza" is $I(D_L)$

$$\begin{aligned} I(D_L) &= 80 \times (1 - (0.375)^2 - (0.625)^2) \\ &= 80 \times (1 - 0.140625 - 0.390625) \\ &= 80 \times (1 - 0.53125) \\ &= 80 \times 0.46875 \end{aligned}$$

$$I(D_L) = 37.5$$

The impurity on "doesn't like pizza" is $I(D_R)$

$$\begin{aligned} I(D_R) &= 20 \times (1 - (0.5)^2 - (0.5)^2) \\ &= 20 \times (1 - 0.25 - 0.25) \\ &= 20 \times (1 - 0.5) \\ &= 20 \times 0.5 = 10 \end{aligned}$$

$$I(D_R) = 10$$

Thus,

$$\begin{aligned} G_1 &= I(D) - (I(D_L) + I(D_R)) \\ &= 48 - 37.5 - 10 \\ &= 48 - \cancel{48} - 47.5 \\ &= 0.5 \end{aligned}$$

So, we should choose the noisy attribute since it has the largest G_1 .

b) a_1 will be the root of the tree and left branch denotes $a_1 = 0$, right branch denotes $a_1 = 1$.

The decision tree which avoids overfitting would have a single decision on the root corresponding to a_1 . This is because a_1 is a attribute predictive of the outcome, where 1% can be considered as noise, and none of the other attributes are predictive of the outcome.

2) a) Let $T = \{t_z | z = 1, \dots, K\}$. Let's assume $\textcircled{5}$ that a data point $P \in S_{ij}$ with corresponding clustering center t_{ij} and $P \in S_2$ with corresponding clustering center t_2 . By triangular Inequality for Euclidean distance, we have.

$$\|P - t_2\|_2 \leq \|P - t_{ij}\|_2 + \|t_{ij} - t_2\|_2$$

Then we square the inequality on both sides,

$$\|P - t_2\|_2^2 \leq (\|P - t_{ij}\|_2 + \|t_{ij} - t_2\|_2)^2 \leq 2\|P - t_{ij}\|_2^2 + 2\|t_{ij} - t_2\|_2^2$$

Then we sum up over all $P \in S_{ij}$, we have

$$\sum_{P \in S_{ij}} \|P - t_2\|_2^2 \leq 2 \sum_{P \in S_{ij}} \|P - t_{ij}\|_2^2 + 2 \sum_{P \in S_{ij}} \|t_{ij} - t_2\|_2^2$$

$$= 2 \sum_{P \in S_{ij}} \|P - t_{ij}\|_2^2 + 2|S_{ij}| \times \|t_{ij} - t_2\|_2^2$$

Summing up j from 1 to K and i from 1 to L ,

$$\sum_{i=1}^L \sum_{j=1}^K \sum_{P \in S_{ij}} \|P - t_2\|_2^2 \leq 2 \sum_{i=1}^L \sum_{j=1}^K \sum_{P \in S_{ij}} \|P - t_{ij}\|_2^2 +$$

$$2 \sum_{i=1}^L \sum_{j=1}^K \omega(t_{ij}) \|t_{ij} - t_2\|_2^2$$

$$\Rightarrow \omega(t_{ij}) = |S_{ij}|$$

The left side of inequality can be reformulated as, (6)

$$\sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|P - t_2\|_2^2 = \sum_{P \in S} \|P - t_2\|_2^2 = \sum_{P \in S} d(P, T).$$

The right side of the inequality can be reformulated as,

$$\begin{aligned} 2 \sum_{i=1}^l \sum_{j=1}^k w(t_{ij}) \|t_{ij} - t_2\|_2^2 + 2 \sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|P - t_{ij}\|_2^2 \\ = 2 \sum_{t_{ij} \in S} w(t_{ij}) d(t_{ij}, T) + 2 \sum_{i=1}^l \sum_{P \in S_i} d(P, T_i) \end{aligned}$$

By cost function,

$$\sum_{P \in S} d(P, T) = \text{cost}(S, T)$$

$$\sum_{t_{ij} \in S} w(t_{ij}) d(t_{ij}, T) = \text{cost}_w(S, T)$$

$$\sum_{i=1}^l \sum_{P \in S_i} d(P, T_i) = \sum_{i=1}^l \text{cost}(S_i, T_i)$$

So, we prove that

$$\text{cost}(S, T) \leq 2 \cdot \text{cost}_w(S, T) + 2 \sum_{i=1}^l \text{cost}(S_i, T_i)$$

2) b) Since ALG is an α -approx ⑦.

$$\text{cost}(S_i, T_i) \leq \alpha \min_{|T'|=k} \{ \text{cost}(S_i, T') \} \leq \alpha \text{cost}(S_i, T_i^*)$$

Then we sum up i from 1 to l

$$\sum_{i=1}^l \text{cost}(S_i, T_i) \leq \sum_{i=1}^l \alpha \cdot \text{cost}(S_i, T_i^*) = \alpha \cdot \text{cost}(S, T^*)$$

2) c) Since we use ALG on \hat{S} and ALG is a α -approx algo

$$\text{cost}_w(\hat{S}, T) \leq \alpha \cdot \min_{|T'|=k} \{ \text{cost}(\hat{S}, T') \} \leq \alpha \cdot \text{cost}_w(\hat{S}, T^*)$$

Hence, first hint proved.

Let $T^* = \{t_2^* \mid z = 1, \dots, k\}$. Suppose a data point $p \in S_{ij}$ with corresponding clustering center t_{ij} and $p \in S_2^*$ with corresponding clustering center t_2^* . By triangular inequality for Euclidean distance,

$$\|t_{ij} - t_2^*\|_2 \leq \|p - t_{ij}\|_2 + \|p - t_2^*\|_2$$

we square the inequalities on both sides, (8)

$$\|t_{ij} - t_2^*\|_2^2 \leq (\|P - t_{ij}\|_2 + \|P - t_2^*\|_2)^2 \leq 2 \cdot \|P - t_{ij}\|_2^2 + 2 \|P - t_2^*\|_2^2$$

Then we sum up all $P \in S_{ij}$, j from 1 to k and i from 1 to l ,

$$\sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|t_{ij} - t_2^*\|_2^2 \leq 2 \cdot \sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|P - t_{ij}\|_2^2 + 2 \sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|P - t_2^*\|_2^2$$

The left side of the inequality can be reformulated,

$$\sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|t_{ij} - t_2^*\|_2^2 = \sum_{i=1}^l \sum_{j=1}^k |S_{ij}| \cdot \|t_{ij} - t_2^*\|_2^2 = \sum_{t_{ij} \in \tilde{S}} w(t_{ij}) \|t_{ij} - t_2^*\|_2^2$$

$$= \sum_{t_{ij} \in \tilde{S}} w(t_{ij}) d(t_{ij}, T^*)$$

The right side of inequality can be reformulated

$$2 \sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|P - t_{ij}\|_2^2 + 2 \sum_{i=1}^l \sum_{j=1}^k \sum_{P \in S_{ij}} \|P - t_2^*\|_2^2 = 2 \sum_{i=1}^l \sum_{P \in S_i} d(P, T_i) + 2 \sum_{P \in S} \phi(P, T^*)$$

By cost function

(9)

$$\sum_{k_{ij} \in S} w(k_{ij}) d(k_{ij}, T^*) = \text{cost}_w(\hat{S}, T^*)$$

$$\sum_{i=1}^l \sum_{p \in S_i} d(p, \hat{T}_i) = \sum_{\hat{T}_i=1}^l \text{cost}(S_i, \hat{T}_i)$$

$$\sum_{p \in S} d(p, T^*) = \text{cost}(S, T^*)$$

So, proving second limit

$$\text{cost}_w(\hat{S}, T^*) \leq 2 \sum_{i=1}^l \text{cost}(S_i, \hat{T}_i) + 2 \cdot \text{cost}(S, T^*)$$

~~By using (a) + (b)~~

$$\begin{aligned} \text{cost}(S, T) &\leq 2 \cdot \text{cost}_w(\hat{S}, T) + 2 \sum_{i=1}^l \text{cost}(S_i, \hat{T}_i) \\ &\leq 2 + \text{cost}_w(\hat{S}, T^*) + 2 \sum_{i=1}^l \text{cost}(S_i, \hat{T}_i) \\ &\leq 2 + \left(2 \sum_{i=1}^l \text{cost}(S_i, \hat{T}_i) + 2 \cdot \text{cost}(S, T^*) \right) + 2 \sum_{i=1}^l \text{cost}(S_i, \hat{T}_i) \\ &\leq 4 + \text{cost}(S, T^*) + 4 + \sum_{i=1}^l \text{cost}(S_i, \hat{T}_i) \\ &\quad + 2 \sum_{i=1}^l \text{cost}(S_i, \hat{T}_i) \end{aligned}$$

$$\leq (4\delta^2 + 6\delta) \cdot \cos t (S, T^*) \quad (10)$$

Thus,

$$\cos t (S, T) \leq (4\delta^2 + 6\delta) \cdot \cos t (S, T^*)$$

$$\begin{aligned} 3) \ a) \quad P_r(\tilde{F}[i] \leq F[i] + \epsilon t) &= 1 - P_r(\tilde{F}[i] \geq F[i] + \epsilon t) \\ &= 1 - P_r(\min_j \{C_j, h_j(i)\} \geq F[i] + \epsilon t) \\ &= 1 - P_r(C_j, h_j(i) \geq F[i] + \epsilon t, \\ &\quad \forall 1 \leq j \leq \lceil \log(\frac{1}{\delta}) \rceil) \\ &= 1 - \prod_{j=1}^{\lceil \log(\frac{1}{\delta}) \rceil} P_r(C_j, h_j(i) \geq F[i] + \epsilon t). \end{aligned}$$

The last equality is followed by independence of hash functions.

By Markov's inequality

$$P_r(C_j, h_j(i) \geq F[i] + \epsilon t) \leq \frac{E[C_j, h_j(i) - F[i]]}{\epsilon t}$$

By second property,

$$E[C_j h_j(i)] \leq F[i] + \frac{\epsilon}{e} (t - F[i]) \quad (11)$$

$$\Leftrightarrow E[C_j h_j(i) - F[i]] \leq \frac{\epsilon}{e} (t - F[i]) < \frac{\epsilon t}{e}$$

Thus,

$$P_r(C_j h_j(i) \geq F[i] + \epsilon t) \leq \frac{E[C_j h_j(i) - F[i]]}{\epsilon t} \leq \frac{1}{e}$$

Insert $P_r(C_j h_j(i) \geq F[i] + \epsilon t) \leq \frac{1}{e}$ into 1st equation

$$P_r(\tilde{F}[i] \leq F[i] + \epsilon t) = 1 - \prod_{j=1}^{\lceil \log(\frac{1}{\delta}) \rceil} P_r(C_j h_j(i) \geq F[i] + \epsilon t)$$

$$= 1 - \left(\frac{1}{e} \right)^{\lceil \log(\frac{1}{\delta}) \rceil} \geq 1 - \delta$$

we prove that

$$P_r(\tilde{F}[i] \leq F[i] + \epsilon t) \geq 1 - \delta$$

3) b)

12

Relative error vs exact word frequency (n = 5, n_buckets = 10000)

