

1.

(a) Show that maximizing the divergence of the label distribution in the child nodes to the uniform distribution is equivalent to minimizing entropy.

Sol:

P: true distribution data

Q: model, approximation of P

$$\text{Entropy}(S) = -\sum_{i=1}^k p_i \log_2 p_i$$

$$D_{KL}(P||Q) = \sum_{i=1}^k p_i \log_2 \frac{p_i}{Q} = \sum_{i=1}^k p_i \log_2 p_i - \sum_{i=1}^k p_i \log_2 Q \text{ (constant)}$$

$$= \sum_{i=1}^k p_i \log_2 p_i - \log_2 Q \sum_{i=1}^k p_i$$

$$\Rightarrow \min \text{Entropy} \sim \max D_{KL}(P||Q)$$

(b) Show that maximizing information gain is equivalent minimizing entropy.

S splits to S_L and S_R

$$\begin{aligned} \min_S H(S_L, S_R) &= \min_S \left[\frac{|S_L|}{|S|} H(S_L) + \frac{|S_R|}{|S|} H(S_R) \right] \\ &= \max_A G(S, A) = \max_A \left[H(S) - \sum_{j=1}^k P(A=j) H(S_j) \right] \\ &= \max_A \left[H(S) - P(A=i) H(S_i) - P(A \neq j) H(S_j) \right] \\ &= \max_A \left[H(S) - \frac{|S_A|}{|S|} H(S_i) - \frac{|S_{A'}|}{|S|} H(S_j) \right] \\ &= H(S) - \min_A \frac{|S_A|}{|S|} H(S_i) \end{aligned}$$

(c) What is the time complexity to find the best split using entropy assuming binary features? / continuous features?

$O(nd)$ for binary classification

$O(n^2 d)$ for continuous regression

2(a)

$$\begin{aligned} L(D) &= \frac{1}{|D|} \sum_{(\vec{x}_j, y_j) \in D} (y_j - \frac{1}{|S|} \sum_{(\vec{x}_i, y_i) \in D} y_i)^2 \\ &= \frac{1}{|D|} \sum_{(\vec{x}_j, y_j) \in D} [y_j^2 - 2y_j \frac{1}{|S|} \sum_{(\vec{x}_i, y_i) \in D} y_i + (\frac{1}{|S|} \sum_{(\vec{x}_i, y_i) \in D} y_i)^2] \\ L(S) &= \frac{1}{|S|} \sum_{(\vec{x}_i, y_i) \in D} (y_i - p)^2 \end{aligned}$$

$$L'(S) = \frac{1}{|S|} \sum_{(\bar{x}_L, y_i) \in D} (y_i - p) \quad (-1) = 0$$

$$p = \frac{\sum_{(\bar{x}_L, y_i) \in D} y_i}{|S|} \quad (\text{average}) \quad \text{and} \quad L''(S) > 0$$

$$L(S): \text{convex}, p = \frac{\sum_{(\bar{x}_L, y_i) \in D} y_i}{|S|}, \quad L(S) \text{ achieves global minimum at the point}$$

2 (b)

(i)

L: x_1, x_2, \dots, x_i

R: $x_{i+1}, x_{i+2}, \dots, x_n$

$$\bar{y}_L = \frac{1}{i} \sum_{j=1}^i y_j$$

$$\bar{y}_R = \frac{1}{n-i} \sum_{j=i+1}^n y_j$$

(ii)

$$\mathcal{L}_L^i: \frac{1}{i} \sum_{j=1}^i (y_j - \bar{y}_L)^2$$

$$\mathcal{L}_R^i: \frac{1}{n-i} \sum_{j=i+1}^n (y_j - \bar{y}_R)^2$$

O(n - i)

(iii)

$$\bar{y}_L^{i+1} = \frac{1}{i+1} \sum_{j=1}^{i+1} y_j = \frac{1}{i+1} \sum_{j=1}^i y_j + y_{i+1} = \frac{1}{i+1} (i\bar{y}_L^i + y_{i+1})$$

(iv)

$$\mathcal{L}_L^{i+1} = \frac{1}{i+1} \sum_{j=1}^{i+1} (y_j - \bar{y}_L^{i+1})^2 = \frac{1}{i+1} \sum_{j=1}^{i+1} (y_j^2 - 2y_j\bar{y}_L^{i+1} + (\bar{y}_L^{i+1})^2)$$

$$s^i = \sum_{j=1}^i y_j^2$$

$$y_i \bar{y}_L^{i+1} = \frac{1}{i+1} (s^i + y_{i+1}^2) - \frac{2y_L^{i+1}}{i+1} \sum_{j=1}^{i+1} y_j$$

$$= s^i - 2\bar{y}_L^{i+1} \bar{y}_L^{i+1} + (\bar{y}_L^{i+1})^2$$

$$= s^i - (\bar{y}_L^{i+1})^2$$

(v)

$$s_{i+1} = \sum_{j=1}^{i+1} y_j^2 = s^i + y_{i+1}^2$$

For each c_i : $O(1) \rightarrow O(nd)$

$$O(n) + O(1)O(n) = O(n)$$

$$O(d \ n \log n) + O(dn) = O(d \ n \log n) \text{ better than } O(dn^2)$$