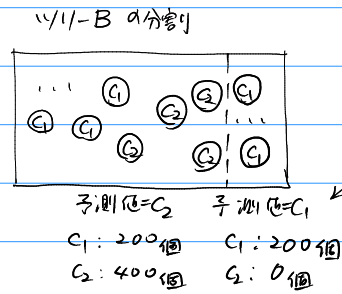
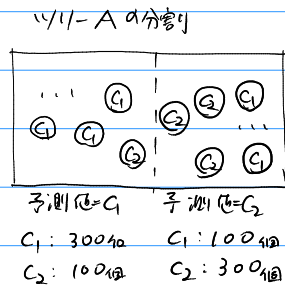


14.11



予測値 = 領域内で最も多い方

11-A の左の葉ノードの誤分類率は

$$E_L^A = \frac{1}{N_{RL}^A} \sum_{\alpha \in ER_L^A} I(\tau_\alpha \neq C_1) = \frac{100}{400} = 0.25$$

11-A の右の葉ノードの誤分類率は

$$E_R^A = \frac{1}{N_{RR}^A} \sum_{\alpha \in ER_R^A} I(\tau_\alpha \neq C_2) = \frac{100}{400} = 0.25$$

11-B の左の葉ノードの誤分類率は

$$E_L^B = \frac{1}{N_{RL}^B} \sum_{\alpha \in ER_L^B} I(\tau_\alpha \neq C_2) = \frac{200}{600} = 0.33$$

11-B の右の葉ノードの誤分類率は

$$E_R^B = \frac{1}{N_{RR}^B} \sum_{\alpha \in ER_R^B} I(\tau_\alpha \neq C_1) = \frac{0}{200} = 0$$

← 問題文には反する

とある。11-A, B で各ノードの誤分類率が特に等しいというわけではない。

(11/1-Aに717)

左右の葉1-1の点割合

$$p_{LC_1}^A = \frac{1}{N_L^A} \sum_{x_i \in R_L^A} I(x_i = C_1) = \frac{300}{400} = 0.75$$

$$p_{LC_2}^A = \text{上と同様} = \frac{100}{400} = 0.25$$

$$p_{RC_1}^A = \text{上と同様} = \frac{100}{400} = 0.25$$

$$p_{RC_2}^A = \text{上と同様} = \frac{300}{400} = 0.75$$

交差エントロピー-損失関数は

$$Q_L^{\text{Across}} = - \sum_{k \in \{C_1, C_2\}} p_{Lk}^A \ln p_{Lk}^A = - (p_{LC_1}^A \ln p_{LC_1}^A + p_{LC_2}^A \ln p_{LC_2}^A) \\ = - (0.75 \ln 0.75 + 0.25 \ln 0.25) = 0.56$$

$$Q_R^{\text{Across}} = - (p_{RC_1}^A \ln p_{RC_1}^A + p_{RC_2}^A \ln p_{RC_2}^A) \\ = - (0.25 \ln 0.25 + 0.75 \ln 0.75) = 0.56$$

このとき極小化基準は (14.31)式は N_C 個の項を213

$$C^{\text{Across}}(T) = \sum_{C=1}^{|T|} N_C Q_C(T) + \lambda |T| = 400 \cdot Q_L^{\text{Across}} + 400 \cdot Q_R^{\text{Across}} + \lambda \cdot 2 = 448 + 2\lambda$$

この係数は

$$Q_L^{\text{AGini}} = \sum_{k \in \{C_1, C_2\}} p_{Lk}^A (1 - p_{Lk}^A) = p_{LC_1}^A (1 - p_{LC_1}^A) + p_{LC_2}^A (1 - p_{LC_2}^A) \\ = 0.75 \times 0.25 + 0.25 \times 0.75 = 0.38$$

$$Q_R^{\text{AGini}} = p_{RC_1}^A (1 - p_{RC_1}^A) + p_{RC_2}^A (1 - p_{RC_2}^A) = 0.25 \times 0.75 + 0.75 \times 0.25 = 0.38$$

このとき極小化基準は

$$C^{\text{AGini}}(T) = 400 \cdot Q_L^{\text{AGini}} + 400 \cdot Q_R^{\text{AGini}} + \lambda \cdot 2 = 304 + 2\lambda$$

(1111-Bに7117)

左右の葉1-1のデータ点割合

$$p_{LC_1}^B = \frac{1}{N_L^B} \sum_{x_i \in R^B} I(x_i = C_1) = \frac{200}{600} = 0.33$$

$$p_{LC_2}^B = \text{上と同様} = \frac{400}{600} = 0.67$$

$$p_{RC_1}^B = \text{上と同様} = \frac{200}{200} = 1$$

$$p_{RC_2}^B = \text{上と同様} = \frac{0}{200} = 0$$

交差エントロピー損失関数は

$$Q_L^{\text{cross}} = - \sum_{k \in \{C_1, C_2\}} p_{Lk}^B \ln p_{Lk}^B = - (p_{LC_1}^B \ln p_{LC_1}^B + p_{LC_2}^B \ln p_{LC_2}^B) \\ = - (0.33 \ln 0.33 + 0.67 \ln 0.67) = 0.63$$

$$Q_R^{\text{cross}} = - (p_{RC_1}^B \ln p_{RC_1}^B + p_{RC_2}^B \ln p_{RC_2}^B) \\ = - (1 \ln 1 + 0 \ln 0) = 0$$

このとき極小化基準は

$$C^{\text{cross}}(T) = \sum_{t \in T} N_t Q_t(T) + \lambda |T| = 600 \cdot 0.63 + 200 \cdot 0 + \lambda \cdot 2 = 378 + 2\lambda$$

ジニ係数は

$$Q_L^{\text{Gini}} = \sum_{k \in \{C_1, C_2\}} p_{Lk}^B (1 - p_{Lk}^B) = p_{LC_1}^B (1 - p_{LC_1}^B) + p_{LC_2}^B (1 - p_{LC_2}^B) \\ = 0.33 \cdot 0.67 + 0.67 \cdot 0.33 = 0.44$$

$$Q_R^{\text{Gini}} = p_{RC_1}^B (1 - p_{RC_1}^B) + p_{RC_2}^B (1 - p_{RC_2}^B) = 0$$

このとき極小化基準は

$$C^{\text{Gini}}(T) = \sum_{t \in T} N_t Q_t(T) + \lambda |T| = 600 \cdot 0.44 + 200 \cdot 0 + \lambda \cdot 2 = 264 + 2\lambda$$

とすると

変数 λ を用いた枝刈り基準は

$$\text{ツリー-A において } 448 + 2\lambda$$

$$\text{ツリー-B において } 378 + 2\lambda$$

したがって ツリー-B のほうが小さくなっている

この係数を用いた枝刈り基準は

$$\text{ツリー-A において } 304 + 2\lambda$$

$$\text{ツリー-B において } 264 + 2\lambda$$

したがって ツリー-B のほうが小さくなっている