

AA1 695 - HW-

Q.1 Bias-variance tradeoff

Bias is the inability of a machine learning model to capture the true relationship between data points due to its inherent assumptions.

eg: A linear regression model assumes a linear relationship between variables

Variance is the difference in how a model fits over new test data & training data, an indicator of overfitting.

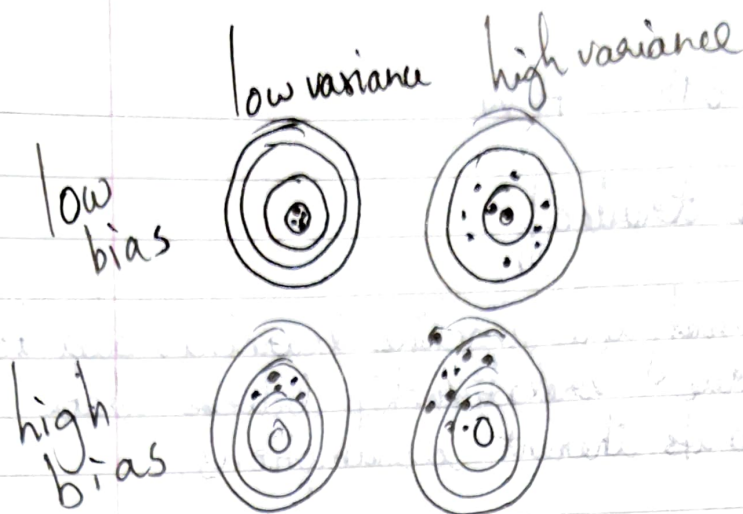
High bias: assuming more about the relⁿ b/w variables

Low bias: fewer assumptions by the learning algo

High bias + low variance: underfitting
(can't capture the relⁿ b/w variables)

Low bias + high variance: overfitting
(can't generalize over new data)

We need to find a good balance between the bias and variance of the model, without overfitting or underfitting, this is called the bias-variance tradeoff.



- to fix bias, adding more i/p's will help the model fit better.
- adding more polynomials can make it more complex to the model
- to fix variance, reducing i/p features or using features with more importance to reduce overfitting
- more training data will also help.

Q.2

actual	50	30
	40	60
predicted		

$$\text{precision} = \frac{TP}{TP+FP} = \frac{50}{40+50} = 0.555$$

$$\text{recall} = \frac{TP}{TP+FN} = \frac{50}{50+30} = 0.625$$

$$f_1\text{-score} = 2 \times (0.555 \times 0.625) / (0.625 + 0.555) = 0.587$$

Ans. 3 target \rightarrow play (6P, 4N)

$$\text{entropy} = \frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$\text{info gain} = \text{entropy}(S) - \sum_{r \in \text{values}(A) | S|} \frac{|S_v|}{|S|} \text{entropy}(S_v)$$

$$\begin{aligned} \text{entropy of play} &= \frac{6}{10} \log_2 \frac{6}{10} - \frac{4}{10} \log_2 \frac{4}{10} \\ &= 0.973 \end{aligned}$$

$$\begin{aligned} \text{entropy of sunny} &= \frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} \\ (1P, 3N) &= 0.817 \end{aligned}$$

$$\text{entropy of overcast} = \frac{2}{2} \log_2 1 - 0 = 0$$

$$\begin{aligned} \text{entropy of rain} &= \frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \\ &= 0.817 \end{aligned}$$

$$\text{gain}(S, \text{outlook}) = E_S - \sum_{|S_v|} \frac{|S_v|}{|S|} \text{entropy } S_v$$

$$= 0.973 - \frac{4}{10} \times 0.817 - \frac{4}{10} \times 0.817$$

$$= \boxed{0.3196}$$

$$E_{\text{Hot}}(1P, 2N) = \frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.93$$

$$E_{\text{mid}} (2P, 1N) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= 0.93$$

$$E_{\text{cool}} (2P, 1N) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.93$$

$$E_{\text{high}} (2P, 3N) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= \underline{\underline{0.973}}$$

$$E_{\text{Normal}} (4P, 1N) = -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5}$$

$$= 0.72$$

$$\text{gain}(s, \text{temp}) = 0.973 - \frac{3}{10} \times 0.93 \times 3 = \boxed{0.136}$$

$$\text{gain}(s, \text{humidity}) = 0.973 - \frac{5}{10} \times 0.973 - \frac{5}{10} \times 0.72$$

$$E_{\text{weak}} (5P, 2N) = \boxed{0.126}$$

$$= 0.87$$

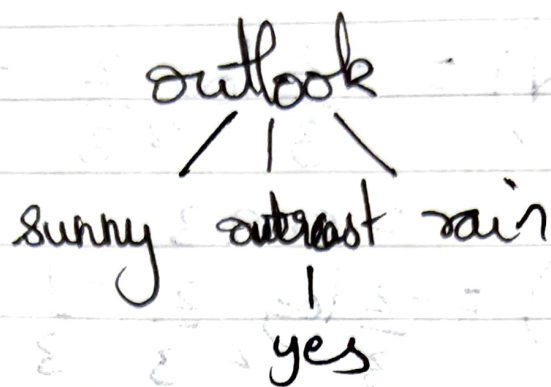
$$E_{\text{strong}} (1P, 2N) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.93$$

$$\text{gain}(s, \text{wind}) = 0.973 - \frac{7}{10} \times 0.876$$

$$= \boxed{0.081} - \frac{3}{10} \times 0.93$$

gain (s, outlook) has highest gain



sunny, humidity:

$$E_{\text{high}}(OP, 3N) = 0, E_{\text{normal}}(IP, 0N) = 0$$

$$\text{gain}(\text{sunny, humid}) = 0.817$$

$$E_{\text{Hot}}(OP, 2N) = 0, E_{\text{mid}}(OP, 1N) = 0, E_{\text{cool}}(IP, 0N) = 0$$

$$\text{gain}(\text{sunny, temp}) = 0.817$$

$$E_{\text{weak}}(IP, 2N) = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3}$$

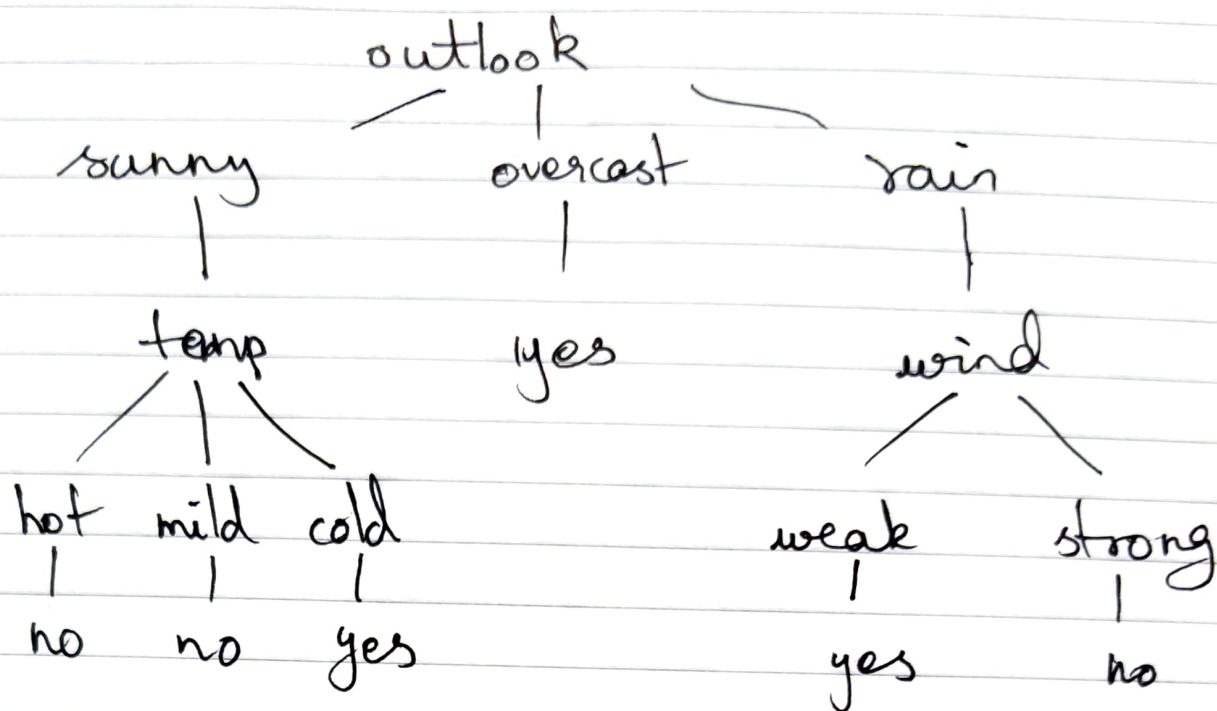
$$= 0.93$$

$$E_{\text{strong}}(3P, 1N) = 0$$

$$\text{gain}(\text{sunny, windy}) = 0.817 - \frac{3}{4} \times 0.93$$

$$= 0.1195$$

gain(sunny, humid) + gain(sunny, temp) will have
info gain



Ans 4 $P(w_1 | d_{11}(x)) = \frac{40}{40+30} = 0.57$

$$P(w_1 | d_{21}(x)) = \frac{20}{40} = 0.5$$

$$P(w_1 | d_{32}(x)) = 0/10 = 0$$

$$w(\text{class 1}) = 0.57 \times 0.5 \times 0 = 0$$

$$P(w_2 | d_{11}(x)) = \frac{30}{30+40} = 0.428$$

$$P(w_2 | d_{21}(x)) = \frac{20}{30+40} = 0.5$$

$$P(w_2 | d_{32}(x)) = \frac{10}{10+0} = 1$$

$$w(\text{class 2}) = 0.48 \times 0.5 \times 1 = 0.24$$

class 2 has higher preference