(a). Run :

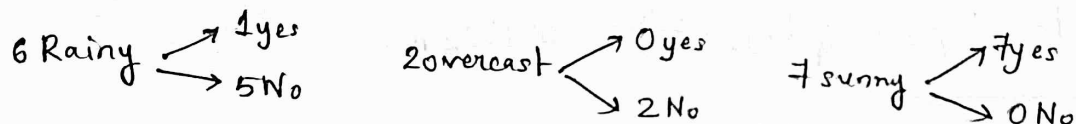| | Yes | No | Total |
|---|---|---|---|
| | 8 | 7 | 15 |

$$E(s) = \sum p(x) \log_2 \frac{1}{p(x)}$$

$$E(s) = -\frac{7}{15} \log_2 \frac{7}{15} - \frac{8}{15} \log_2 \frac{8}{15} = 0.99679$$

$$IG(s, outlook) = H(s) - \sum p(x) * H(x) \qquad \boxed{H(s) = 0.99679}$$

| Out = Rainy | Out = Overcast | Out = Sunny | Total |
|---|---|---|---|
| 6 | 2 | 7 | 15 |

$$P(s_{rainy}) = \frac{6}{15} \qquad P(s_{overcast}) = \frac{2}{15} \qquad P(s_{sunny}) = \frac{7}{15}$$

6 Rainy ⟶ 1 yes / 5 No

2 overcast ⟶ 0 yes / 2 No

7 sunny ⟶ 7 yes / 0 No

$$H(s_{rainy}) = -\frac{5}{6} \log_2 \left(\frac{5}{6}\right) - \frac{1}{6} \log_2 \left(\frac{1}{6}\right) = 0.6500$$
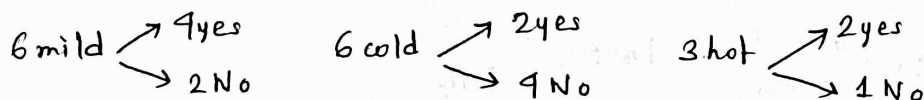
$$H(s_{overcast}) = -\log(1) - 0 = 0$$

$$H(s_{sunny}) = -\log(1) - 0 = 0$$

$$IG(s, outlook) = 0.99679 - \frac{6}{15} \times 0.65 \qquad \therefore \boxed{IG(s, outlook) = 0.73679}$$

$$IG(s, temp) = H(s) - \sum p(x) * H(x) \qquad H(s) = 0.99679$$

| temp = mild | temp = cold | temp = hot | total |
|---|---|---|---|
| 6 | 6 | 3 | 15 |

$$P(s_{mild}) = \frac{6}{15} \qquad P(s_{cold}) = \frac{6}{15} \qquad P(s_{temp}) = \frac{3}{15}$$

6 mild ⟶ 4 yes / 2 No

6 cold ⟶ 2 yes / 4 No

3 hot ⟶ 2 yes / 1 No

$$H(s_{mild}) = -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) = 0.918296$$

$$H(s_{cold}) = -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) = 0.918296$$

$$H(s_{hot}) = -\frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right) = 0.918296$$
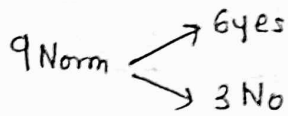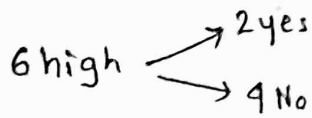
$$IG(s, temp) = 0.99679 - \frac{12}{15}(0.918296) - \frac{3}{15}(0.918296) = \boxed{IG(s, temp) = 0.07899}$$

$IG(s, hum) = H(s) - \sum p(x) * H(x)$

$H(s) = 0.99679$.

| hum= high | hum = norm | total |
|---|---|---|
| 6 | 9 | 15 |

$P(S_{high}) = \frac{6}{15}$    $P(S_{norm}) = \frac{9}{15}$

6 high $\nearrow$ 2 yes $\searrow$ 4 No     9 Norm $\nearrow$ 6 yes $\searrow$ 3 No

$H(S_{high}) = -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) = 0.918296$

$H(S_{nom}) = -\frac{6}{9} \log_2 \left(\frac{6}{9}\right) - \frac{3}{9} \log_2 \left(\frac{3}{9}\right) = 0.918296$.

$\therefore$ $\boxed{IG(s, Norm) = 0.07849}$.

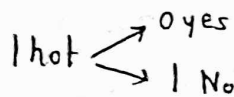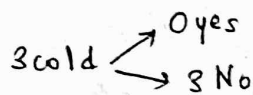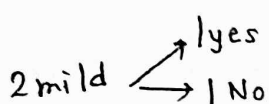$\therefore$ the first node will be outlook. (highest info gain).

| Rainy : | Yes | No | Total |
|---|---|---|---|
| | 1 | 5 | 6 |

$H(S_{rainy}) = -\frac{5}{6} \log_2 \left(\frac{5}{6}\right) - \frac{1}{6} \log_2 \left(\frac{1}{6}\right) = 0.65$    $\boxed{H(S_{rainy}) = 0.65}$.

$IG(S_{rainy}, temp) = H(S_{rainy}) - \sum p(x) * H(x)$.

| temp = mild | temp = cold | temp = hot | total |
|---|---|---|---|
| 2 | 3 | 1 | 6 |

$P(S_{mild}) = \frac{2}{6}$    $P(S_{cold}) = \frac{3}{6}$    $P(S_{hot}) = \frac{1}{6}$

2 mild $\nearrow$ 1 yes $\searrow$ 1 No    3 cold $\nearrow$ 0 yes $\searrow$ 3 No    1 hot $\nearrow$ 0 yes $\searrow$ 1 No

$H(S_{mild}) = -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) = 1$

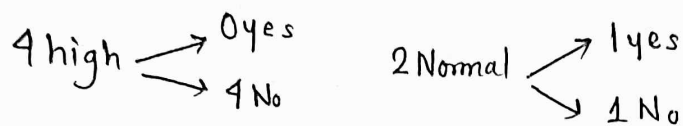$H(S_{cold}) = -\log 1 - 0 = 0$

$H(S_{hot}) = -\log 1 - 0 = 0$

$IG(S_{rainy}, temp) = 0.65 - \left(1 \times \frac{2}{6}\right) = 0.3166$    $\boxed{IG(S_{rainy}, temp) = 0.3166}$.

$IG(S_{rainy}, hum) = H(S_{rainy}) - \sum p(x) * H(x)$

$H(S_{rainy}) = 0.65$

| hum = high | hum = normal | total |
|---|---|---|
| 4 | 2 | 6 |

$P(S_{high}) = \frac{4}{6}$ $\qquad$ $P(S_{nom}) = \frac{2}{6}$

4 high $\begin{cases} \to \text{0yes} \\ \to \text{4 No} \end{cases}$ $\qquad$ 2 Normal $\begin{cases} \to \text{1yes} \\ \to \text{1 No} \end{cases}$
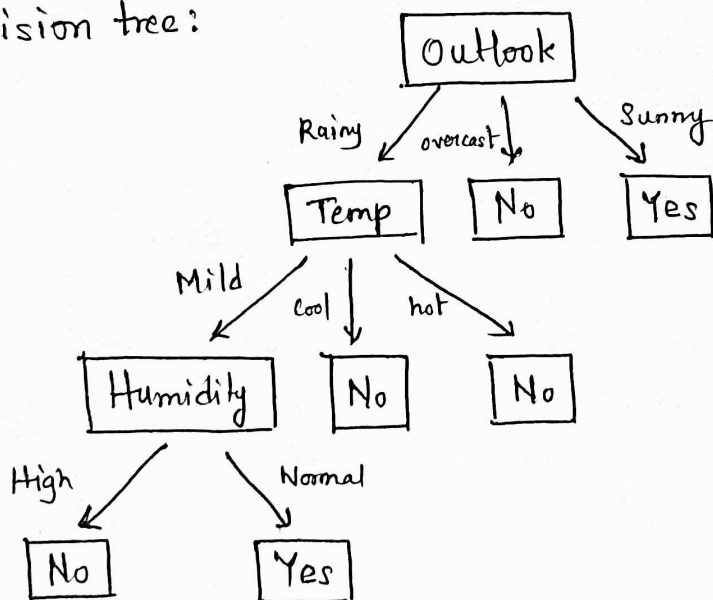
$H(S_{high}) = \log 1 - 0 = 0$

$H(S_{nomal}) = -\frac{1}{2}\log\frac{1}{2} - \frac{1}{2}\log\frac{1}{2} = 1$

$IG(S_{rainy}, nomal) = 0.65 - (1 \times \frac{2}{6}) = 0.3166$ $\qquad$ $\boxed{IG(S_{rainy}, hum) = 0.3166}$

Therefore, the next decision node can be any of humidity or temperature.

∴ The decision tree:



(b) If tomorrows, outlook is overcast, the temperature is cool and the Humidity level is high, the runner will not go for a run.

2.(a)

**Accuracy -**

Training/validation accuracy for minimum node entropy 0.010000 is 1.000 / 0.863

Training/validation accuracy for minimum node entropy 0.050000 is 0.999 / 0.863

Training/validation accuracy for minimum node entropy 0.100000 is 0.997 / 0.865

Training/validation accuracy for minimum node entropy 0.200000 is 0.990 / 0.867

Training/validation accuracy for minimum node entropy 0.400000 is 0.979 / 0.861

Training/validation accuracy for minimum node entropy 0.800000 is 0.919 / 0.856

Training/validation accuracy for minimum node entropy 1.000000 is 0.871 / 0.840

Training/validation accuracy for minimum node entropy 2.000000 is 0.596 / 0.600

Test accuracy with minimum node entropy 0.200000 is 0.872

**Error Rate –**

Training/validation error rate for minimum node entropy 0.010000 is 0.000 / 0.137

Training/validation error rate for minimum node entropy 0.050000 is 0.001 / 0.137

Training/validation error rate for minimum node entropy 0.100000 is 0.003 / 0.135

Training/validation error rate for minimum node entropy 0.200000 is 0.010 / 0.133

Training/validation error rate for minimum node entropy 0.400000 is 0.021 / 0.139

Training/validation error rate for minimum node entropy 0.800000 is 0.081 / 0.144

Training/validation error rate for minimum node entropy 1.000000 is 0.129 / 0.160

Training/validation error rate for minimum node entropy 2.000000 is 0.404 / 0.400
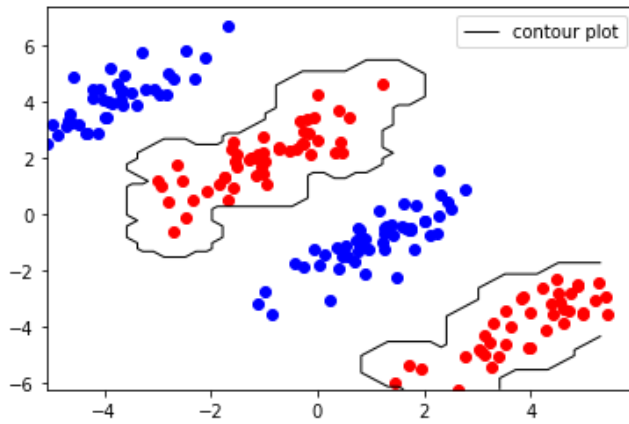
Test error rate with minimum node entropy 0.200000 is 0.128

2.(b) With the increase in the minimum node entropy, the validation accuracy decreases for both validation and training dataset. If there was a high error-rate on both training and validation data that would indicate that the model may be too simple. But that does not happen in the above case. If the error-rate is slightly increased on training data compared to validation data, thus the model is not too complex as it seems.

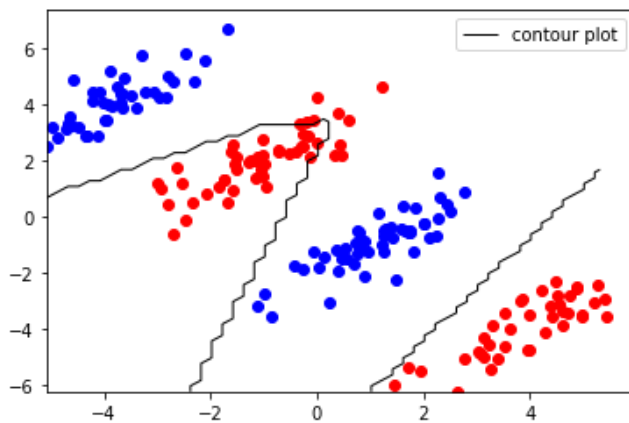3.(a) sigma = 0.001

Accuracy on the simulated data is 1.00

Decision Boundary –
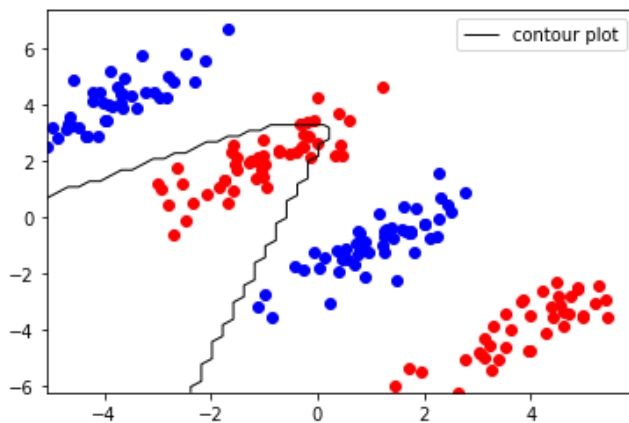


sigma = 0.1

Accuracy on the simulated data is 0.95

Decision Boundary –
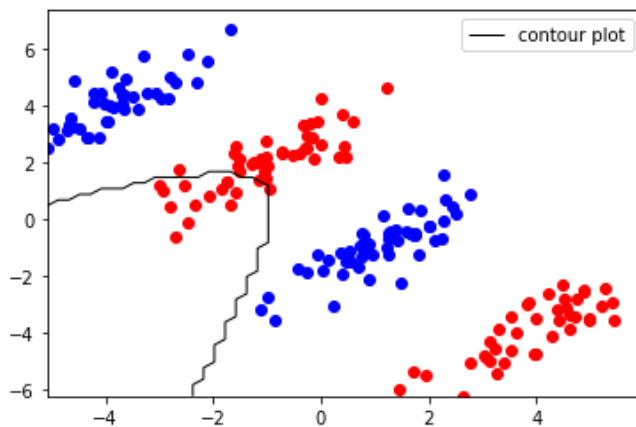


sigma = 1

Accuracy on the simulated data is 0.72

Decision Boundary –

sigma = 10

Accuracy on the simulated data is 0.57

Decision Boundary –



3.(b) Accuracy for sigma pool= (0.001,0.001,0.005) – 0.51 (on second dataset containing 7 and 9)

Accuracy for sigma pool = (0.1,0.1,0.5) – 0.83 (on second dataset containing 7 and 9)

Accuracy for sigma pool= (1,1,5) – 0.96 (on second dataset containing 7 and 9)

Accuracy for sigma pool= (10,1,5) – 0.96 (on second dataset containing 7 and 9)

With increasing sigma, the decision boundaries are not accurate as can be seen from the above contour plots. The accuracy of the model on the simulated data decreases as the sigma is increased. Thus, we have to optimize sigma so that the validation accuracy increases again it does not overfit to the training data set.

Whereas on decreasing sigma the accuracy on the training set does not increase but rather decreases (for the second dataset containing 7 and 9) whereas it does not seem to change on the first dataset (containing 4 and 9) which remains at a low value throughout.

3.(c) sigma pool= (0.001,0.001,0.005):

Accuracy on the digit classification-49 is 0.50

Accuracy on the digit classification-79 is 0.51

sigma pool = (0.1,0.1,0.5):

Accuracy on the digit classification-49 is 0.50

Accuracy on the digit classification-79 is 0.83

sigma pool= (1,1,5):

Accuracy on the digit classification-49 is 0.50

Accuracy on the digit classification-79 is 0.96

sigma pool= (10,1,5):

Accuracy on the digit classification-49 is 0.5

Accuracy on the digit classification-79 is 0.96