

1) Decision stump

$$h(x) = \begin{cases} + & \text{if sneezing} = \text{Yes} \\ - & \text{if sneezing} = \text{No} \end{cases}$$

Dataset:

1. (Yes, +)  $\rightarrow$  predicted = +  $\rightarrow$  correct
2. (No, -)  $\rightarrow$  predicted = -  $\rightarrow$  correct
3. (Yes, -)  $\rightarrow$  predicted = +  $\rightarrow$  wrong
4. (No, +)  $\rightarrow$  predicted = -  $\rightarrow$  correct

Step-1: Count Errors

total sample = 4

misclassified sample = 1 (sneezed 3 only)

step-2: Training error rate.

$$\begin{aligned} \text{Train [Error rate]} &= \frac{\text{Train \{ \#errors \}}}{\text{Train \{ total sample \}}} \\ &= \frac{1}{4} = 0.25 = 25\% \end{aligned}$$

Step-3:- Compare with memorizer model

Memorizer model: remembers all training data  $\rightarrow$  predicts perfectly  $\rightarrow$  0% error.

Decision stump: 25% Error

a) Training error rate = 25%

b) Memorizer is better (0% error).

Q2).

split on Age (x1)

Groups:-

young : records {1,2}  $\rightarrow$  labels = {yes, yes}

Majority = yes  $\rightarrow$  Errors in this group = 0

Mid : records {3,6}  $\rightarrow$  labels = {No, No}

Majority = No  $\rightarrow$  Errors = 0

old : records {4,5}  $\rightarrow$  labels = {No, Yes}

\* If tie / majority break by majority  $\rightarrow$  we must pick one label:

majority counts: 1 No, 1 Yes  $\rightarrow$  tie.

\* In practice for training - error split, choose the majority label (tie implies any choice causes 1 error) with either choice you'll get 1 misclassified in this group.

\* Errors = 1

$$\text{Total Errors} = 0 + 0 + 1 = 1$$

$$\text{Training Error rate} = 1/6 \approx 0.1667 \approx 16.67\%$$

Split on Exercise (x2)

Groups:

High : records {1,5}  $\rightarrow$  labels = {yes, yes}  $\rightarrow$  Majority, Yes  $\rightarrow$  Errors = 0

Medium : records {2,4}  $\rightarrow$  labels = {yes, No}  $\rightarrow$  Majority?

tie  $\rightarrow$  whichever label chosen produces 1 error  $\rightarrow$  Error = 1

Low : Records {3,6}  $\rightarrow$  labels = {No, No}  $\rightarrow$  Majority No  $\rightarrow$  Errors = 0

$$\text{Total Errors} = 0 + 1 + 0 = 1$$

$$\text{Training error rate} = 1/6 = 16.67\%$$



split on diet (x3)

Groups:-

Poor : records  $\{1, 3, 4, 6\} \rightarrow$  labels =  $\{No, No, No, No\}$

Counts : Yes = 1, No = 3  $\rightarrow$  Majority = No  $\rightarrow$  Errors = #Yes = 1

Good : records  $\{2, 5\} \rightarrow$  labels =  $\{Yes, Yes\} \rightarrow$  Majority = Yes  $\rightarrow$  Errors = 0

Total errors = 1 + 0 = 1.

1) Training error rates for splitting on each feature

split on Age  $\Rightarrow 1/6 = 16.67\%$ .

split on Exercise  $\rightarrow 1/6 = 16.67\%$ .

split on Diet (x3)  $\rightarrow 1/6 = 16.67\%$ .

2) All the three features tie with the same training error (16.67%). So there is no single best root split by the training-error criterion - any of them is equally good under this metric.

### Q3) Entropy & Information Gain

→ labels: 3 yes, 3 no →

$$H(Y) = -[0.5 \log_2 0.5 + 0.5 \log_2 0.5] = 1.0$$

→ Split on Exercise (x2)

High: (Yes, Yes) → Entropy 0

Medium: (Yes, No) → Entropy 1.0

Low: (No, No) → Entropy 0

$$\text{weighted Entropy} = (2/6)0 + (2/6)1 + (2/6)0 = 1/3 = 0.333$$

$$\rightarrow \text{Information gain} = 1 - 0.333 = 0.667$$

Exercise is a good split

### Q4) Confusion Matrix metrics:

Confusion matrix ( $T_P=25$ ,  $F_N=5$ ,  $F_P=15$ ,  $T_N=55$ , total=100)

$$\text{Accuracy} = (25 + 55) / 100 = 0.80$$

$$\text{Precision} = 25 / (25 + 15) = 0.625$$

$$\text{Recall} = 25 / (25 + 5) = 0.833$$

$$\text{specificity} = 55 / (55 + 15) = 0.786$$

$$F_1 = 2 \cdot (0.625 \cdot 0.833) / (0.625 + 0.833) = 0.714$$

If imbalanced (80 negatives, 20 positives) Recall, Precision &  $F_1$  are more informative than accuracy.



### Q5) Distance calculations (KNN)

New point  $P(5,4)$

$$d(P,A) = \sqrt{(5-2)^2 + (4-4)^2} = \sqrt{9} = 3$$

$$d(P,B) = \sqrt{(5-4)^2 + (4-4)^2} = \sqrt{1} = 1$$

$$d(P,C) = \sqrt{(5-4)^2 + (4-6)^2} = \sqrt{5} = 2.236$$

1-NN = nearest neighbor is B (blue)  $\rightarrow$  Predict blue

3-NN = neighbours = {Red, blue, Red}  $\rightarrow$  Majority Red

### Q6) K-fold Cross-Validation

Average errors:

$$K=1 \rightarrow (0.20 + 0.25 + 0.15 + 0.30) / 4 = 0.225$$

$$K=3 \rightarrow (0.15 + 0.20 + 0.10 + 0.20) / 4 = 0.1625$$

$$K=5 \rightarrow (0.10 + 0.15 + 0.10 + 0.20) / 4 = 0.1375$$

Best Generalization ~~k~~  $K=5$ .