

Group Members:

1- Mert Tazeoğlu - b21946606

Q8

a) In find-S algorithm, we ignore negative samples. Therefore will deal with 1st, 2nd and 3rd samples.

$$h_0 = \{ \langle \emptyset, \emptyset, \emptyset, \emptyset \rangle, \langle \emptyset, \emptyset, \emptyset, \emptyset \rangle \}$$

$$S_1 = \{ \langle \text{mole}, \text{brown}, \text{tall}, \text{US} \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \}$$

$$h_1 = \{ \langle \text{mole}, \text{brown}, \text{tall}, \text{US} \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \}$$

$$S_2 = \{ \langle \text{mole}, \text{brown}, \text{short}, \text{French} \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \}$$

$$h_2 = \{ \langle \text{mole}, \text{brown}, ?, ? \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \}$$

$$S_3 = \{ \langle \text{ferret}, \text{brown}, \text{tall}, \text{German} \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{indian} \rangle \} \Rightarrow \text{since it is negative we ignore it}$$

$$h_3 = \{ \langle \text{mole}, \text{brown}, ?, ? \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \} \Rightarrow \text{therefore it is same with } h_2$$

$$S_4 = \{ \langle \text{mole}, \text{brown}, \text{tall}, \text{Irish} \rangle, \langle \text{ferret}, \text{brown}, \text{short}, \text{Irish} \rangle \}$$

$$h_4 = \{ \langle \text{mole}, \text{brown}, ?, ? \rangle, \langle \text{ferret}, ?, \text{short}, ? \rangle \}$$

$$b) S_0 = \{ \langle \emptyset, \emptyset, \emptyset, \emptyset \rangle, \langle \emptyset, \emptyset, \emptyset, \emptyset \rangle \}$$

$$G_0 = \{ \langle ?, ?, ?, ? \rangle, \langle ?, ?, ?, ? \rangle \}$$

$$S_1 = \{ \langle \text{mole}, \text{brown}, \text{tall}, \text{US} \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \}$$

$$G_1 = \{ \langle \langle ?, ?, ?, ? \rangle, \langle ?, ?, ?, ? \rangle \}$$

$$S_2 = \{ \langle \text{mole}, \text{brown}, ?, ? \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \}$$

$$G_2 = \{ \langle \langle ?, ?, ?, ? \rangle, \langle ?, ?, ?, ? \rangle \}$$

$$S_3 = \{ \langle \text{mole}, \text{brown}, ?, ? \rangle, \langle \text{ferret}, \text{black}, \text{short}, \text{US} \rangle \}$$

$$G_3 = \{ \langle \langle \text{mole}, ?, ?, ? \rangle, \langle ?, ?, ?, ? \rangle \rangle, \langle \langle ?, ?, ?, ? \rangle, \langle ?, ?, ?, ? \rangle \rangle \}$$

$$S_4 = \{ \langle \text{mole}, \text{brown}, ?, ? \rangle, \langle \text{ferret}, ?, \text{short}, ? \rangle \}$$

$$G_4 = \{ \langle \langle \text{mole}, ?, ?, ? \rangle, \langle ?, ?, ?, ? \rangle \}$$

$$S_4 \times G_4 = \{ \langle \langle \text{mole}, ?, ?, ? \rangle, \langle ?, ?, ?, ? \rangle \rangle \}$$

$$c) 2^8 = 256$$

Each hypothesis one consistent with that example can have either specified value seen (given) or "?" for each attribute. therefore for each of 8 attributes, can take 2 different values. therefore $2^8 = 256$ distinct hypotheses are consistent with that.

Q2

a) Label is same of point x with 1-NN algorithm. Since 1-NN chooses only one nearest point, then there can't be another nearer point to x other than both of nearest neighbours of $S1$ and $S2$.

Let $S1$'s nearest neighbours as:

$x \rightarrow (+)$, distance = 1

$y \rightarrow (-)$, distance = 2

$z \rightarrow (-)$, distance = 5

then, 1-NN($S1$) is positive.

Let $S2$'s nearest neighbours as:

$x \rightarrow (+)$, distance = 1

$t \rightarrow (-)$, distance = 2

$q \rightarrow (-)$, distance = 3

then, 1-NN($S2$) is positive.

therefore, $(S1 \cup S2)$'s nearest neighbours are:

$x \rightarrow (+)$, distance = 1

$y \rightarrow (-)$, distance = 2

$t \rightarrow (-)$, distance = 2

} then 1-NN($S1 \cup S2$) is positive and it will consider only point x .

b) Let $S1$'s nearest neighbours are:

$a \rightarrow (-)$, distance = 1

$b \rightarrow (-)$, distance = 2

$c \rightarrow (+)$, distance = 4

then, 3-NN($S1$) = positive.

Let $S2$'s nearest neighbours are:

$d \rightarrow (-)$, distance = 3

$e \rightarrow (+)$, distance = 5

$f \rightarrow (+)$, distance = 6

then, 3-NN($S2$) = positive.

therefore, $(S1 \cup S2)$'s nearest neighbours are:

$a \rightarrow (-)$, distance = 1

$b \rightarrow (-)$, distance = 2

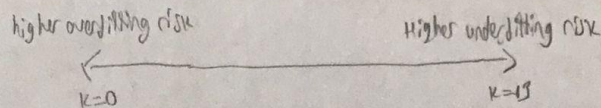
$d \rightarrow (-)$, distance = 3

} then 3-NN($S1 \cup S2$) is negative.

Q3

a) A point can be its own neighbor, therefore $k=0$ minimizes training set error.
 If $k=0$, then training error is 0.

b) Too big k values (for example $k=13$) misclassifies every datapoint (using leave one out cross validation). On the other hand, too small k value (for example $k=0$) leads to overfitting.



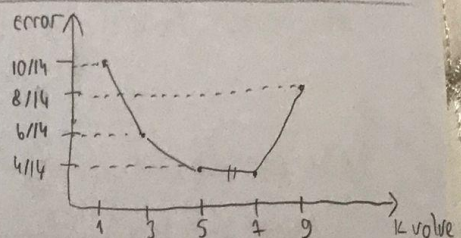
c) for $k=1$ error is $10/14$

for $k=3$ error is $6/14$

for $k=5$ error is $4/14$

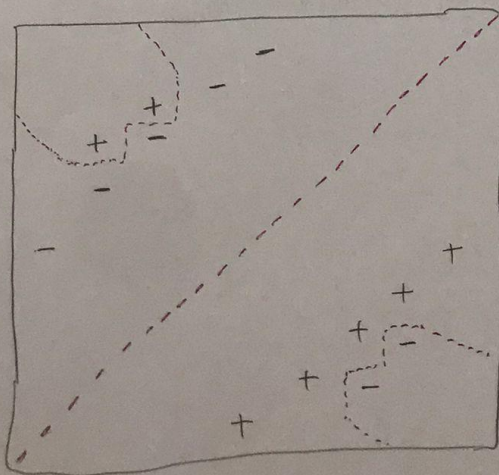
for $k=7$ error is $4/14$

for $k=9$ error is $8/14$



* $k=5$ and $k=7$ minimizes the leave-one-out cross validation error. Error is $4/14$.

d)



$\text{---} \Rightarrow 1$ nearest neighbour decision boundary

(Q4)

	ACTUAL	
	positive	negative
PREDICTED		
positive	TP	FP
negative	FN	TN

In 1000 instances,
 i) Since 950 are healthy, TP+FN must be 950.
 ii) Since 50 are patient, FP+TN must be 50.

a)

	+	-
+	TP = 950	FP = 50
-	FN = 0	TN = 0

In this example;

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN} = \frac{950}{1000} = 95\%$$

According to accuracy value, this model is very good. However, it is quite bad since it classifies every person as positive. (In other words, it doesn't detect negative)

b)

	+	-
+	TP = 50	FP = 1
-	FN = 900	TN = 49

In this example;

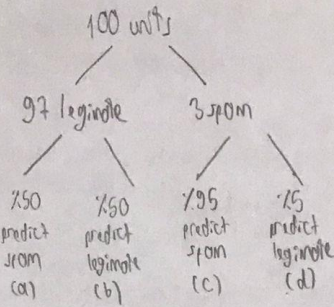
$$\text{Recall} = \frac{TP}{TP + FN} = \frac{50}{950} = 0.052 \text{ (so bad value)}$$

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{50}{51} = 0.98 \text{ (good value)}$$

In this example, model has a good precision value, but a bad recall value. In this example, using only precision value not enough and both of them should be used. If we use only precision values, although the model is bad it will be evaluated as good.

Q5

Solution-1:



$$a = (97 \times 50) / 100 = \%48.5$$

$$b = (97 \times 50) / 100 = \%48.5$$

$$c = (3 \times 95) / 100 = \%2.85$$

$$d = (3 \times 5) / 100 = \%0.15$$

$$\text{Predict As Spam} \Rightarrow a + c = \%51.35$$

$$\text{Actually spam (predicted as spam)} \Rightarrow \%2.85$$

$$\frac{2.85}{51.35} = \%5.55$$

Solution-2:

This can be solved using Bayes Theorem, which says:

$P(A|B) = P(B|A) \times P(A) / P(B)$, where A = email is spam and B = algorithm predict as spam

$$P(A|B) = P(B|A) \times P(A) / P(B)$$

$$P(B|A) = 95/100 = 0.95$$

$$P(A) = 3/100 = 0.03$$

$$P(B) = 3/100$$

$$P(B) = \frac{(0.97 \times 0.5) + (0.03 \times 0.95)}{1 \times 1.00} = \frac{0.4850 + 0.0285}{1} = 0.5135$$

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} = \frac{(0.95) \times (0.03)}{(0.5135)} = \frac{0.0285}{0.5135} \approx 0.555 \approx \%5.55$$

Q6 (PART-1)

Recursion depth: 1 (root)

Grouped records: {Age: {young: 5, middle: 4, senior: 5}, sex: {M: 4, F: 4}, BP: {high: 4, normal: 6, low: 4}, cholesterol: {normal: 8, high: 6}}

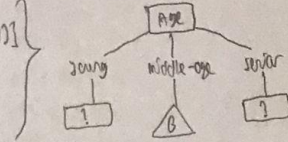
System entropy: 0.918

Feature entropies: {Age: {0.99, 0.91, 0.91}, sex: {0.98, 0.99}, BP: {1.0, 0.65, 0.81}, cholesterol: {0.81, 0.91}}

Information gains: {Age: 0.95, sex: 0.4, BP: 0.6, cholesterol: 0.005}

Most informative: Age (young > senior > middle)

↳ in middle age drug B is common



Recursion depth: 2 (young people)

Grouped records: {sex: {F: 3, M: 1}, BP: {high: 1, normal: 2, low: 0}, cholesterol: {normal: 2, high: 2}}

System entropy: 0.918

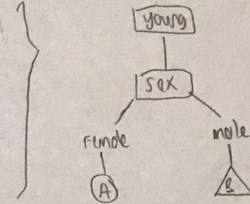
Feature entropies: {sex: {0.0, 0.0}, BP: {0.0, 0.0}, cholesterol: {0.0, 0.0}}

Information gains: {sex: 0.918, BP: 0.44, cholesterol: 0.44}

Most informative: sex (female > male)

↳ All of young males use drug B.

Also all of young females use drug A.



Recursion depth: 2 (senior people)

Grouped records: {sex: {F: 2, M: 1}, BP: {normal: 3, low: 1}, cholesterol: {normal: 3, high: 1}}

System entropy: 0.864

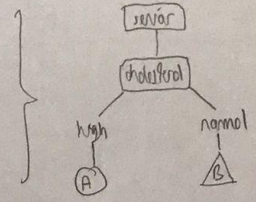
Feature entropies: {sex: {0.0, 0.0}, BP: {0.0, 0.0}, cholesterol: {0.0, 0.0}}

Information gains: {sex: 0.19, BP: 0.44, cholesterol: 0.44}

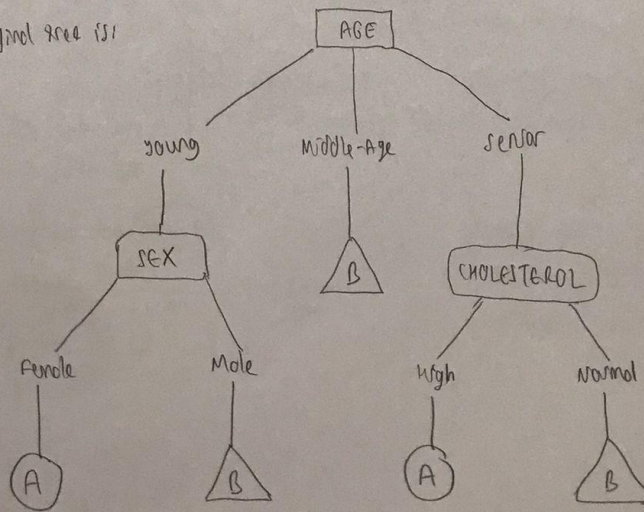
Most informative: cholesterol (high > low)

↳ there is a exact match that normal cholesterol = drug B, otherwise drug A

no exact match in BP but cholesterol provides it therefore we select it



At the end, our final tree is:



(P15) ⇒ since P15 is middle aged, then he/she should use DRUG-B!

$$\begin{array}{l}
 b) \quad P(\text{drug} = A) = 5/14 \\
 P(\text{drug} = B) = 9/14
 \end{array}
 \quad
 \begin{array}{l}
 P(\text{cholesterol} = \text{normal}) = 8/14 \\
 P(\text{cholesterol} = \text{high}) = 6/14
 \end{array}
 \quad
 \begin{array}{l}
 P(\text{BP} = \text{low}) = 5/14 \\
 P(\text{BP} = \text{normal}) = 4/14 \\
 P(\text{BP} = \text{high}) = 5/14
 \end{array}
 \quad
 \begin{array}{l}
 P(\text{sex} = F) = 4/14 \\
 P(\text{sex} = M) = 7/14
 \end{array}
 \quad
 \begin{array}{l}
 P(\text{age} = \text{young}) = 5/14 \\
 P(\text{age} = \text{middle}) = 4/14 \\
 P(\text{age} = \text{senior}) = 5/14
 \end{array}$$

Let's compute $P(X_i | c_i)$ for each class:

- $P(\text{age} = \text{'middle'} | \text{drug} = \text{'A'}) = 0/4 = 0$
- $P(\text{age} = \text{'middle'} | \text{drug} = \text{'B'}) = 4/4 = 1$
- $P(\text{sex} = \text{'female'} | \text{drug} = \text{'A'}) = 4/4$
- $P(\text{sex} = \text{'female'} | \text{drug} = \text{'B'}) = 3/7$
- $P(\text{BP} = \text{'low'} | \text{drug} = \text{'A'}) = 1/4$
- $P(\text{BP} = \text{'low'} | \text{drug} = \text{'B'}) = 3/4$
- $P(\text{cholesterol} = \text{'normal'} | \text{drug} = \text{'A'}) = 2/8$
- $P(\text{cholesterol} = \text{'normal'} | \text{drug} = \text{'B'}) = 6/8$

$$X = (\text{age} = \text{'middle'}, \text{sex} = \text{'female'}, \text{BP} = \text{'low'}, \text{cholesterol} = \text{'normal'})$$

$$P(X | c_i): P(X | \text{drug} = \text{'A'}) = 0 \times 4/4 \times 1/4 \times 2/8 = 0.000$$

$$P(X | \text{drug} = \text{'B'}) = 1 \times 3/4 \times 3/4 \times 6/8 = 0.281$$

$$P(X | c_i) \times P(c_i) \div P(X | \text{drug} = \text{'A'}) \times P(\text{drug} = \text{'A'}) = 0.000$$

$$P(X | \text{drug} = \text{'B'}) \times P(\text{drug} = \text{'B'}) = 0.154$$

→ therefore P15 belongs to class drug = B.

$$\text{confidence of classification} = \frac{0.154}{0.154 + 0.000} = 1 = \underline{\underline{100\%}}$$

BOTH OF MODELS

SAY THAT P15

SHOULD USE PRUG-B!