

--	--	--	--	--	--	--	--	--	--	--

**Internal Assessment Test –III, January 2022**

Sub:	MACHINE LEARNING						Code:	18MCA53	
Date:	24-01-2022	Duration:	90 mins	Max Marks:	50	Sem:	V	Branch:	MCA

**Answer Any 5 Questions**

	Marks	OBE	
		CO	RBT
<b>Part-I</b>			
Q1 Illustrate Bayesian Belief network with the help of an example.	10	CO4	L2
or			
Q2 Explain Expectation Maximization algorithm.	10	CO4	L2
<b>Part-II</b>			
Q3 Write short note on: a) Sample Error   b) True Error   c) Variance   d) Expected value   e) Confidence Interval	10	CO5	L2
or			
Q4 What is instance-based learning? Explain Case-Based reasoning method.	10	CO5	L1
<b>Part-III</b>			
Q5 Discuss the method of comparing two algorithms. Justify with paired-t method.	10	CO5	L3
or			
Q6 Discuss Q-learning algorithm with the help of an example.	10	CO5	L3
<b>Part-IV</b>			
Q7 Explain the Distance Weighted Nearest Neighbor Algorithm.	10	CO5	L2
or			
Q8 Explain the locally weighted linear regression.	10	CO5	L2
<b>Part-V</b>			
Q9 What is reinforcement learning algorithm? Explain.	10	CO5	L2
or			
Q10 Explain the Radial basis function (RBF) network.	10	CO5	L3

--	--	--	--	--	--	--	--	--	--	--	--

Internal Assessment Test 3 – December. 2020

<b>Sub:</b>	MACHINE LEARNING							<b>Sub Code:</b>	18MCA 53
<b>Date:</b>	15- 12- 20 20	<b>Duration:</b>	90 min's	<b>Max Marks:</b>	50	<b>Sem</b>	5 <sup>th</sup>	<b>Branch:</b>	MCA

**Note : Answer FIVE FULL Questions, choosing ONE full question from each Module**

PART I		MAR KS	OBE	
			C O	RB T
1)	<ul style="list-style-type: none"> <li>Consider the setting in which we wish to learn a nondeterministic (probabilistic) function <math>f : X \rightarrow \{0, 1\}</math>, which has two discrete output values.</li> <li>We want a function approximator whose output is the probability that <math>f(x) = 1</math>. In other words, learn the target function <math>f^* : X \rightarrow [0, 1]</math> such that <math>f^*(x) = P(f(x) = 1)</math></li> </ul> <p><i>How can we learn <math>f^*</math> using a neural network?</i></p> <ul style="list-style-type: none"> <li>Use of brute force way would be to first collect the observed frequencies of 1's and 0's for each possible value of <math>x</math> and to then train the neural network to output the target frequency for each <math>x</math>.</li> </ul> <p><i>What criterion should we optimize in order to find a maximum likelihood hypothesis for <math>f^*</math> in this setting?</i></p> <ul style="list-style-type: none"> <li>First obtain an expression for <math>P(D h)</math></li> <li>Assume the training data <math>D</math> is of the form <math>D = \{(x_1, d_1) \dots (x_m, d_m)\}</math>, where <math>d_i</math> is the observed 0 or 1 value for <math>f(x_i)</math>.</li> <li>Both <math>x_i</math> and <math>d_i</math> as random variables, and assuming that each training example is drawn independently, we can write <math>P(D h)</math> as</li> </ul> $P(D   h) = \prod_{i=1}^m P(x_i, d_i   h) \quad \text{equ (1)}$ <p>Applying the product rule</p> $P(D   h) = \prod_{i=1}^m P(d_i   h, x_i) P(x_i) \quad \text{equ (2)}$ <p>The probability <math>P(d_i h, x_i)</math></p> $P(d_i h, x_i) = \begin{cases} h(x_i) & \text{if } d_i = 1 \\ (1 - h(x_i)) & \text{if } d_i = 0 \end{cases} \quad \text{equ (3)}$ <p>Re-express it in a more mathematically manipulable form, as</p> $P(d_i h, x_i) = h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} \quad \text{equ (4)}$	10	CO4	L5

Equation (4) to substitute for  $P(d_i | h, x_i)$  in Equation (5) to obtain

$$P(D|h) = \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} P(x_i) \quad \text{equ (5)}$$

We write an expression for the maximum likelihood hypothesis

$$h_{ML} = \operatorname{argmax}_{h \in H} \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} P(x_i)$$

The last term is a constant independent of  $h$ , so it can be dropped

$$h_{ML} = \operatorname{argmax}_{h \in H} \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} \quad \text{equ (6)}$$

It easier to work with the log of the likelihood, yielding

$$h_{ML} = \operatorname{argmax}_{h \in H} \sum_{i=1}^m d_i \ln h(x_i) + (1 - d_i) \ln(1 - h(x_i)) \quad \text{equ (7)}$$

Equation (7) describes the quantity that must be maximized in order to obtain the maximum likelihood hypothesis in our current problem setting.

2) **Step 1: Convert the data into a frequency table.**

Play	Frequency
Yes	9
No	5

Outlook	Yes	No
Sunny	2	3
Overcast	4	0
Rain	3	2

Temperature	Yes	No
Hot	2	2
Mild	4	2
Cool	3	1

Humidity	Yes	No
High	3	4
Normal	6	1

Wind	Yes	No
Strong	3	3
Weak	6	2

**Step 2: Create Likelihood table**

Play	Frequency	Likelihood
Yes	9	9/14
No	5	5/14

Outlook	Yes	No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rain	3/9	2/5

Temperature	Yes	No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

Humidity	Yes	No
High	3/9	4/5
Normal	6/9	1/5

Wind	Yes	No
Strong	3/9	3/5
Weak	6/9	2/5

**Step 3:**

**New instance**

**X= (Outlook = Sunny, Temperature = Cool, Humidity = High, Wind = Strong)**

**Play Tennis = ?**

$$\begin{aligned}
 P(X | \text{Play} = \text{Yes}) &= P(\text{Play} = \text{Yes}) * P(\text{Outlook} = \text{Sunny} | \text{Yes}) * P(\text{Temperature} = \text{Cool} | \text{Yes}) * P(\text{Humidity} = \text{High} | \text{Yes}) * P(\text{Wind} = \text{Strong} | \text{Yes}) \\
 &= 9/14 * 2/9 * 3/9 * 3/9 * 3/9 \\
 &= 0.0053
 \end{aligned}$$

$$\begin{aligned}
 P(X | \text{Play} = \text{No}) &= P(\text{Play} = \text{No}) * P(\text{Outlook} = \text{Sunny} | \text{No}) * P(\text{Temperature} = \text{Cool} | \text{No}) * P(\text{Humidity} = \text{High} | \text{No}) * P(\text{Wind} = \text{Strong} | \text{No}) \\
 &= 5/14 * 3/5 * 1/5 * 4/5 * 3/5 \\
 &= 0.0206
 \end{aligned}$$

So: 0.0206 > 0.0053

Result : X : PlayTennis = No

3)

- The naive Bayes classifier makes significant use of the assumption that the values of the attributes  $a_1 \dots a_n$  are conditionally independent given the target value  $v$ .
- This assumption dramatically reduces the complexity of learning the target function

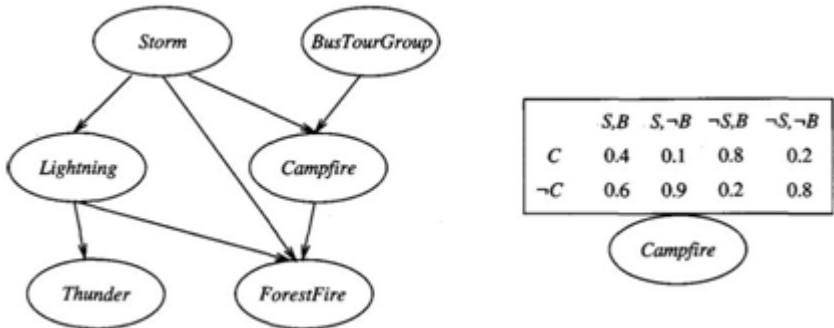
A Bayesian belief network describes the probability distribution governing a set of variables by specifying a set of conditional independence assumptions along with a set of conditional probabilities  
 Bayesian belief networks allow stating conditional independence assumptions that apply to subsets of the variables

**Notation**

- Consider an arbitrary set of random variables  $Y_1 \dots Y_n$ , where each variable  $Y_i$  can take on the set of possible values  $V(Y_i)$ .
- The joint space of the set of variables  $Y$  to be the cross product  $V(Y_1) \times V(Y_2) \times \dots \times V(Y_n)$ .
- In other words, each item in the joint space corresponds to one of the possible assignments of values to the tuple of variables  $(Y_1 \dots Y_n)$ . The probability distribution over this joint' space is called the joint probability distribution.
- The joint probability distribution specifies the probability for each of the possible variable bindings for the tuple  $(Y_1 \dots Y_n)$ .
- A Bayesian belief network describes the joint probability distribution for a set of variables.

**Representation**

A Bayesian belief network represents the joint probability distribution for a set of variables.  
 Bayesian networks (BN) are represented by directed acyclic graphs.



The Bayesian network in above figure represents the joint probability distribution over the boolean variables *Storm, Lightning, Thunder, ForestFire, Campfire, and BusTourGroup*

A Bayesian network (BN) represents the joint probability distribution by specifying a set of *conditional independence assumptions*

- BN represented by a directed acyclic graph, together with sets of local conditional probabilities

- Each variable in the joint space is represented by a node in the Bayesian network
- The network arcs represent the assertion that the variable is conditionally independent of its non-descendants in the network given its immediate predecessors in the network.

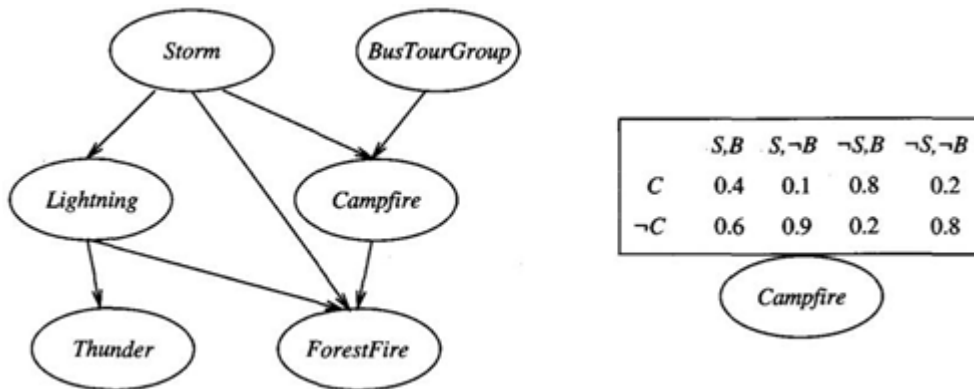
A **conditional probability table (CPT)** is given for each variable, describing the probability distribution for that variable given the values of its immediate predecessors. The joint probability for any desired assignment of values  $(y_1, \dots, y_n)$  to the tuple of network variables  $(Y_1 \dots Y_m)$  can be computed by the formula

$$P(y_1, \dots, y_n) = \prod_{i=1}^n P(y_i | Parents(Y_i))$$

Where,  $Parents(Y_i)$  denotes the set of immediate predecessors of  $Y_i$  in the network.

**Example:**

Consider the node **Campfire**. The network nodes and arcs represent the assertion that **Campfire** is conditionally independent of its non-descendants **Lightning** and **Thunder**, given its immediate parents **Storm** and **BusTourGroup**.



This means that once we know the value of the variables **Storm** and **BusTourGroup**, the variables **Lightning** and **Thunder** provide no additional information about **Campfire**. The conditional probability table associated with the variable **Campfire**. The assertion is

$$P(\text{Campfire} = \text{True} \mid \text{Storm} = \text{True}, \text{BusTourGroup} = \text{True}) = 0.4$$

**Inference**

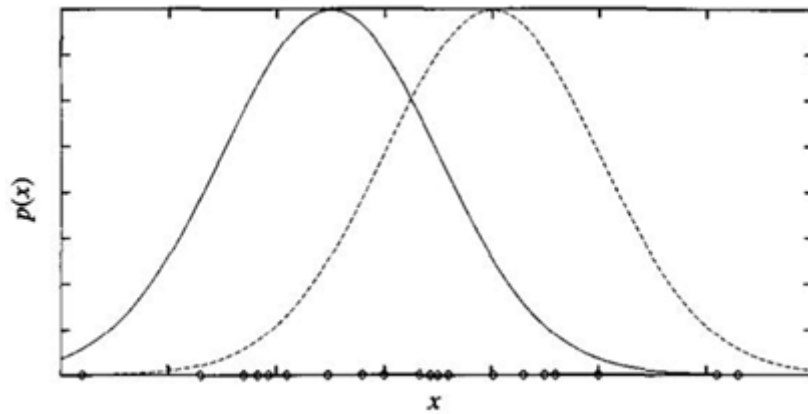
- Use a Bayesian network to infer the value of some target variable (e.g., ForestFire) given the observed values of the other variables.
- Inference can be straightforward if values for all of the other variables in the network are known exactly.
- A Bayesian network can be used to compute the probability distribution for any subset of network variables given the values or distributions for any subset of the remaining variables.
- An arbitrary Bayesian network is known to be NP-hard.

4) **THE EM ALGORITHM**

The EM algorithm can be used even for variables whose value is never directly observed, provided the general form of the probability distribution governing these variables is known.

**Estimating Means of k Gaussians**

Consider a problem in which the data  $D$  is a set of instances generated by a probability distribution that is a mixture of  $k$  distinct Normal distributions.



- This problem setting is illustrated in Figure for the case where  $k = 2$  and where the instances are the points shown along the  $x$  axis.
- Each instance is generated using a two-step process.
  - First, one of the  $k$  Normal distributions is selected at random.
  - Second, a single random instance  $x_i$  is generated according to this selected distribution.
- This process is repeated to generate a set of data points as shown in the figure.
- To simplify, consider the special case
  - The selection of the single Normal distribution at each step is based on choosing each with uniform probability
  - Each of the  $k$  Normal distributions has the same variance  $\sigma^2$ , known value.
- The learning task is to output a hypothesis  $h = (\mu_1, \dots, \mu_k)$  that describes the means of each of the  $k$  distributions.
- We would like to find a maximum likelihood hypothesis for these means; that is, a hypothesis  $h$  that maximizes  $p(D | h)$ .

$$\mu_{ML} = \operatorname{argmin}_{\mu} \sum_{i=1}^m (x_i - \mu)^2 \quad (1)$$

In this case, the sum of squared errors is minimized by the sample mean.

$$\mu_{ML} = \frac{1}{m} \sum_{i=1}^m x_i \quad (2)$$

- Our problem here, however, involves a mixture of  $k$  different Normal distributions, and we cannot observe which instances were generated by which distribution.
- Consider full description of each instance as the triple  $(x_i, z_{i1}, z_{i2})$ ,
  - where  $x_i$  is the observed value of the  $i$ th instance and
  - where  $z_{i1}$  and  $z_{i2}$  indicate which of the two Normal distributions was used to generate the value  $x_i$
- In particular,  $z_{ij}$  has the value 1 if  $x_i$  was created by the  $j^{\text{th}}$  Normal distribution and 0 otherwise.
- Here  $x_i$  is the observed variable in the description of the instance, and  $z_{i1}$  and  $z_{i2}$  are hidden variables.
- If the values of  $z_{i1}$  and  $z_{i2}$  were observed, we could use following Equation to solve for the means  $\mu_1$  and  $\mu_2$
- Because they are not, we will instead use the EM algorithm.

**EM algorithm**

**Step 1:** Calculate the expected value  $E[z_{ij}]$  of each hidden variable  $z_{ij}$ , assuming the current hypothesis  $h = \langle \mu_1, \mu_2 \rangle$  holds.

**Step 2:** Calculate a new maximum likelihood hypothesis  $h' = \langle \mu'_1, \mu'_2 \rangle$ , assuming the value taken on by each hidden variable  $z_{ij}$  is its expected value  $E[z_{ij}]$  calculated in Step 1. Then replace the hypothesis  $h = \langle \mu_1, \mu_2 \rangle$  by the new hypothesis  $h' = \langle \mu'_1, \mu'_2 \rangle$  and iterate.

Let us examine how both of these steps can be implemented. Step 1 must calculate the expected value of each  $z_{ij}$ . This  $E[z_{ij}]$  is the probability ability that instance  $x_i$  was generated by the  $j$ th Normal distribution.

$$E[z_{ij}] = \frac{p(x = x_i | \mu = \mu_j)}{\sum_{n=1}^2 p(x = x_i | \mu = \mu_n)}$$

$$= \frac{e^{-\frac{1}{2\sigma^2}(x_i - \mu_j)^2}}{\sum_{n=1}^2 e^{-\frac{1}{2\sigma^2}(x_i - \mu_n)^2}}$$

Thus the first step is implemented by substituting the current value of the observed  $x_i$  into the above expression.

In the second step we use the  $E[z_{ij}]$  calculated during Step 1 to calculate a new maximum likelihood hypothesis  $h' = \langle \mu'_1, \mu'_2 \rangle$ . The new maximum likelihood hypothesis in this case is given by

$$\mu_j \leftarrow \frac{\sum_{i=1}^m E[z_{ij}] x_i}{\sum_{i=1}^m E[z_{ij}]}$$

5)

**Sample Error –**

The sample error of a hypothesis with respect to some sample  $S$  of instances drawn from  $X$  is the fraction of  $S$  that it misclassifies.

*Definition:* The sample error ( $errors_S(h)$ ) of hypothesis  $h$  with respect to target function  $f$  and data sample  $S$  is

$$errors_S(h) \equiv \frac{1}{n} \sum_{x \in S} \delta(f(x), h(x))$$

Where  $n$  is the number of examples in  $S$ , and the quantity  $\delta(f(x), h(x))$  is 1 if  $f(x) \neq h(x)$ , and 0 otherwise.

**True Error –**

The true error of a hypothesis is the probability that it will misclassify a single randomly drawn instance from the distribution  $\mathbf{D}$ .

*Definition:* The true error ( $error_{\mathbf{D}}(h)$ ) of hypothesis  $h$  with respect to target function  $f$  and distribution  $\mathbf{D}$ , is the probability that  $h$  will misclassify an instance drawn at random according to  $\mathbf{D}$ .

$$error_{\mathbf{D}}(h) \equiv \Pr_{x \in \mathbf{D}} [f(x) \neq h(x)]$$

The Variance captures how far the random variable is expected to vary from its mean value.

10

CO5

L2

**Definition:** The variance of a random variable Y,  $\text{Var}[Y]$ , is

$$\text{Var}[Y] \equiv E[(Y - E[Y])^2]$$

The variance describes the expected squared error in using a single observation of Y to estimate its mean  $E[Y]$ .

The square root of the variance is called the standard deviation of Y, denoted  $\sigma_Y$

The Mean (expected value) is the average of the values taken on by repeatedly sampling the random variable

**Definition:** Consider a random variable Y that takes on the possible values  $y_1, \dots, y_n$ . The expected value (Mean) of Y,  $E[Y]$ , is

$$E[Y] \equiv \sum_{i=1}^n y_i \Pr(Y = y_i)$$

### Confidence Intervals for Discrete-Valued Hypotheses

Suppose we wish to estimate the true error for some discrete valued hypothesis h, based on its observed sample error over a sample S, where

- The sample S contains n examples drawn independent of one another, and independent of h, according to the probability distribution  $\mathbf{D}$
- $n \geq 30$
- Hypothesis h commits r errors over these n examples (i.e.,  $\text{errors}(h) = r/n$ ).

Under these conditions, statistical theory allows to make the following assertions:

1. Given no other information, the most probable value of  $\text{error}_{\mathbf{D}}(h)$  is  $\text{errors}(h)$
2. With approximately **95% probability**, the true error  $\text{error}_{\mathbf{D}}(h)$  lies in the interval

$$\text{errors}_S(h) \pm 1.96 \sqrt{\frac{\text{errors}_S(h)(1 - \text{errors}_S(h))}{n}}$$

#### Example:

Suppose the data sample S contains  $n = 40$  examples and that hypothesis h commits  $r = 12$  errors over this data.

- The **sample error** is  $\text{errors}(h) = r/n = 12/40 = 0.30$
- Given no other information, **true error** is  $\text{error}_{\mathbf{D}}(h) = \text{errors}(h)$ , i.e.,  $\text{error}_{\mathbf{D}}(h) = 0.30$

With the 95% confidence interval estimate for  $\text{error}_{\mathbf{D}}(h)$ .

$$\text{errors}_S(h) \pm 1.96 \sqrt{\frac{\text{errors}_S(h)(1 - \text{errors}_S(h))}{n}}$$

$$= 0.30 \pm (1.96 * 0.07) \quad = 0.30 \pm 0.14$$

3. A different constant,  **$z_N$** , is used to calculate the **N% confidence interval**. The general expression for approximate N% confidence intervals for  $\text{error}_{\mathbf{D}}(h)$  is

$$\text{errors}_S(h) \pm z_N \sqrt{\frac{\text{errors}_S(h)(1 - \text{errors}_S(h))}{n}}$$

#### Example:

Suppose the data sample S contains  $n = 40$  examples and that hypothesis h commits  $r = 12$  errors over this data.



- The **sample error** is  $errors(h) = r/n = 12/40 = 0.30$
- With the 68% confidence interval estimate for  $error_D(h)$

$$errors_S(h) \pm 1.00 \sqrt{\frac{errors_S(h)(1 - errors_S(h))}{n}}$$

$$= 0.30 \pm (1.00 * 0.07)$$

$$= 0.30 \pm 0.07$$

6) **INTRODUCTION**

- Instance-based learning methods such as nearest neighbor and locally weighted regression are conceptually straightforward approaches to approximating real-valued or discrete-valued target functions.
- Learning in these algorithms consists of simply storing the presented training data. When a new query instance is encountered, a set of similar related instances is retrieved from memory and used to classify the new query instance.
- Instance-based approaches can construct a different approximation to the target function for each distinct query instance that must be classified.

**CASE-BASED REASONING**

- Case-based reasoning (CBR) is a learning paradigm based on lazy learning methods and they classify new query instances by analysing similar instances while ignoring instances that are very different from the query.
- In CBR represent instances are not represented as real-valued points, but instead, they use a **rich symbolic** representation.
- CBR has been applied to problems such as conceptual design of mechanical devices based on a stored library of previous designs, reasoning about new legal cases based on previous rulings, and solving planning and scheduling problems by reusing and combining portions of previous solutions to similar problems.

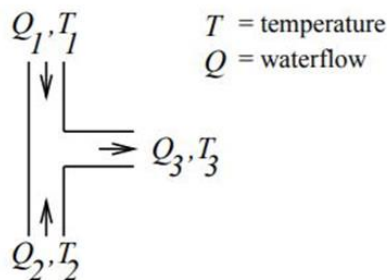
**A prototypical example of a case-based reasoning**

- The CADET system employs case-based reasoning to assist in the conceptual design of simple mechanical devices such as water faucets.
- It uses a library containing approximately 75 previous designs and design fragments to suggest conceptual designs to meet the specifications of new design problems.
- Each instance stored in memory (e.g., a water pipe) is represented by describing both its structure and its qualitative function.
- New design problems are then presented by specifying the desired function and requesting the corresponding structure.

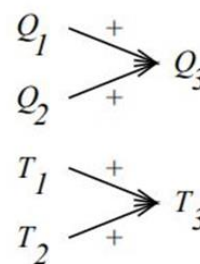
The problem setting is illustrated in below figure

**A stored case: T-junction pipe**

Structure:



Function:



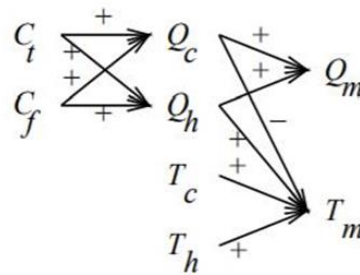
- The function is represented in terms of the qualitative relationships among the water-flow levels and temperatures at its inputs and outputs.
- In the functional description, an arrow with a "+" label indicates that the variable at the arrowhead increases with the variable at its tail. A "-" label indicates that the variable at the head decreases with the variable at the tail.
- Here  $Q_c$  refers to the flow of cold water into the faucet,  $Q_h$  to the input flow of hot water, and  $Q_m$  to the single mixed flow out of the faucet.
- $T_c$ ,  $T_h$ , and  $T_m$  refer to the temperatures of the cold water, hot water, and mixed water respectively.
- The variable  $C_t$  denotes the control signal for temperature that is input to the faucet, and  $C_f$  denotes the control signal for waterflow.
- The controls  $C_t$  and  $C_f$  are to influence the water flows  $Q_c$  and  $Q_h$ , thereby indirectly influencing the faucet output flow  $Q_m$  and temperature  $T_m$ .

**A problem specification: Water faucet**

Structure:

?

Function:



- CADET searches its library for stored cases whose functional descriptions match the design problem. If an exact match is found, indicating that some stored case implements exactly the desired function, then this case can be returned as a suggested solution to the design problem. If no exact match occurs, CADET may find cases that match various subgraphs of the desired functional specification.

7)

**Hypothesis Testing**

- Evaluates 2 mutual exclusive statement on population using sample data.
- Steps:
  1. Make initial Assumptions
  2. Collect Data
  3. Gather Evidence to Reject or accept NULL hypothesis
- What is the probability that  $\text{error}_D(h_1) > \text{error}_D(h_2)$

**COMPARING LEARNING ALGORITHMS**

- Comparing the performance of 2 learning algorithms  $L_A$  and  $L_B$ .
- A reasonable way to define "on average" is to consider the relative performance of these 2 algorithms averaged over all the training sets of size  $n$  over Distribution  $D$ .

$$E_{S \sim D} [\text{error}_D(L_A(S)) - \text{error}_D(L_B(S))]$$

Where ,

$L(S)$  : Hypothesis output of learning method  $L$  when given the sample  $S$  of training data .

Here  $S \sim D$  : The expected value is taken over samples  $S$  drawn according to the underlying instance distribution  $D$ .

10

CO5

L3

## Comparing learning algorithms $L_A$ and $L_B$

1. Partition data  $D_0$  into  $k$  disjoint test sets  $T_1, T_2, \dots, T_k$  of equal size, where this size is at least 30.
2. For  $i$  from 1 to  $k$ , do
  - use  $T_i$  for the test set, and the remaining data for training set  $S_i$
  - $S_i \leftarrow \{D_0 - T_i\}$
  - $h_A \leftarrow L_A(S_i)$
  - $h_B \leftarrow L_B(S_i)$
  - $\delta_i \leftarrow error_{T_i}(h_A) - error_{T_i}(h_B)$
3. Return the value  $\bar{\delta}$ , where

$$\bar{\delta} \equiv \frac{1}{k} \sum_{i=1}^k \delta_i$$

### Paired t Tests

Following is the estimation procedure:

- We are given the observed values of a set of independent, identically distributed random variables  $Y_1, Y_2, \dots, Y_k$ .
- We wish to estimate the mean of the probability distribution governing these  $Y_i$ .
- The estimator will use sample mean which is given by

$$\bar{Y} \equiv \frac{1}{k} \sum_{i=1}^k Y_i$$

- We can request a new training examples drawn according to the underlying instance distribution. Modify the steps on each iteration through the loop it generates a new random training set  $S_i$  and a new random test set  $T_i$  by drawing from this underlying instance distribution instead of drawing sample  $D_0$ .

Confidence interval =

$$\mu = \bar{Y} \pm t_{N,k-1} s_{\bar{Y}}$$

where  $s_{\bar{Y}}$  is the estimated standard deviation of the sample mean

$$s_{\bar{Y}} = \sqrt{\frac{1}{k(k-1)} \sum_{i=1}^k (Y_i - \bar{Y})^2}$$

10)

### The Q Function

The value of Evaluation function  $Q(s, a)$  is the reward received immediately upon executing action  $a$  from state  $s$ , plus the value (discounted by  $\gamma$ ) of following the optimal policy thereafter

$$Q(s, a) \equiv r(s, a) + \gamma V^*(\delta(s, a)) \quad \text{equ (4)}$$

Rewrite Equation (3) in terms of  $Q(s, a)$  as

$$\pi^*(s) = \operatorname{argmax}_a Q(s, a) \quad \text{equ (5)}$$

Equation (5) makes clear, it need only consider each available action  $a$  in its current state  $s$  and choose the action that maximizes  $Q(s, a)$ .

### An Algorithm for Learning Q

- Learning the Q function corresponds to learning the **optimal policy**.

The key problem is finding a reliable way to estimate training values for  $Q$ , given only a sequence of immediate rewards  $r$  spread out over time. This can be accomplished through *iterative approximation*

$$V^*(s) = \max_{a'} Q(s, a')$$

Rewriting Equation

$$Q(s, a) = r(s, a) + \gamma \max_{a'} Q(\delta(s, a), a')$$

### Q learning algorithm

#### Q learning algorithm

For each  $s, a$  initialize the table entry  $\hat{Q}(s, a)$  to zero.

Observe the current state  $s$

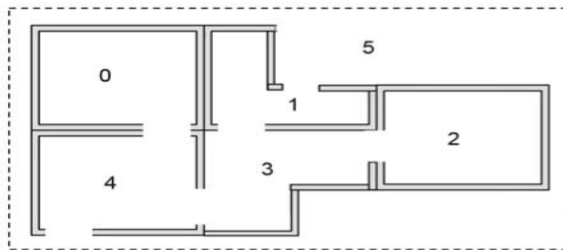
Do forever:

- Select an action  $a$  and execute it
- Receive immediate reward  $r$
- Observe the new state  $s'$
- Update the table entry for  $\hat{Q}(s, a)$  as follows:

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a')$$

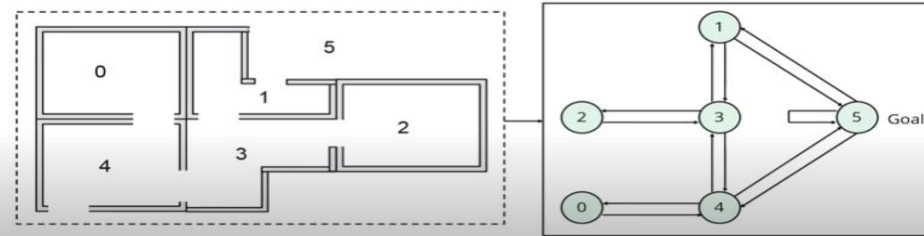
- $s \leftarrow s'$

Place an agent in any one of the rooms (0,1,2,3,4) and the goal is to reach outside the building (room 5)

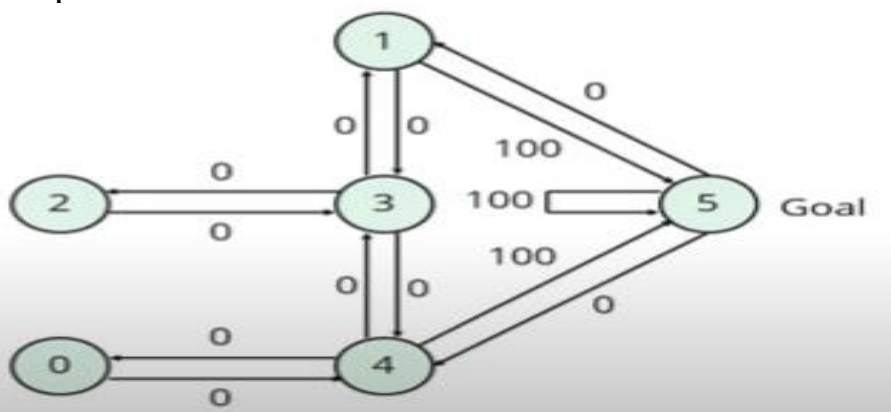


- 5 rooms in a building connected by doors
- each room is numbered 0 through 4
- The outside of the building can be thought of as one big room (5)
- Doors 1 and 4 lead into the building from room 5 (outside)

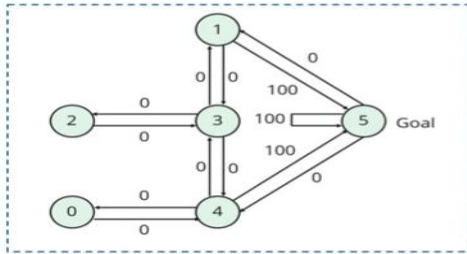
Let's represent the rooms on a graph, each room as a node, and each door as a link



Next Step is to associate reward value to each door.



We can put the state diagram and the instant reward values into a reward table, matrix R.



State	0	1	2	3	4	5
0	-1	-1	-1	-1	0	-1
1	-1	-1	-1	0	-1	100
2	-1	-1	-1	0	-1	-1
3	-1	0	0	-1	0	-1
4	0	-1	-1	0	-1	100
5	-1	0	-1	-1	0	100

The -1's in the table represent null values

**Next Step is to create Q matrix.**

Add another matrix Q, representing the memory of what the agent has learned through experience.

- The rows of matrix Q represent the current state of the agent
- columns represent the possible actions leading to the next state
- Formula to calculate the Q matrix:

$$Q(\text{state}, \text{action}) = R(\text{state}, \text{action}) + \text{Gamma} * \text{Max}[Q(\text{next state}, \text{all actions})]$$

**Note**

The Gamma parameter has a range of 0 to 1 (0 <= Gamma < 1).

- If Gamma is closer to zero, the agent will tend to consider only immediate rewards.
- If Gamma is closer to one, the agent will consider future rewards with greater weight

**Q-Learning Example:**

First step is to set the value of the learning parameter Gamma = 0.8, and the initial state as Room 1.

Next, initialize matrix Q as a zero matrix:

- From room 1 you can either go to room 3 or 5, let's select room 5.
- From room 5, calculate maximum Q value for this next state based on all possible actions:

$$Q(1,5) = R(1,5) + 0.8 * \text{Max}[Q(5,1), Q(5,4), Q(5,5)] = 100 + 0.8 * 0 = 100$$

	0	1	2	3	4	5
0	0	0	0	0	0	0
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
5	0	0	0	0	0	0

State	0	1	2	3	4	5
0	-1	-1	-1	-1	0	-1
1	-1	-1	-1	0	-1	100
2	-1	-1	-1	0	-1	-1
3	-1	0	0	-1	0	-1
4	0	-1	-1	0	-1	100
5	-1	0	-1	-1	0	100

