

Name: - Niloy mallick

ID: - 1731151

CSE 417 - Data mining and warehouse

Final Exam

Q2 → Section A

Supervised vs Unsupervised Learning

Supervised Learning (classification)

— Supervision : the training data (observations, measurements, etc) are accompanied by labels indicating the class of the observations

— New data is classified based on the training set

Unsupervised Learning (clustering)

— The class labels of training data is unknown

— Given a set of measurements, observations, etc with the aim of establishing the existence of classes or clusters in the data.

❑ Evaluating classification methods Issues

Accuracy

- classifier accuracy: predicting class label
- predictor accuracy: guessing values of Predicted attributes.

Speed

- time to construct the model (training time)
- time to use the model (classification / prediction time)

Robustness

- handling noise and missing values

Scalability

- efficiency in disk-resident databases

Interpretability

- understanding and insight provided by the model

Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules.

Q4 → Section B

The Bayes' theorem describes the Probability of an event based on prior knowledge of the conditions that might be relevant to the event. The Bayes' theorem is expressed in the following formula:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Where

* $P(A|B)$ - the Probability of event A occurring, given event B has occurred

* $P(B|A)$ - the Probability of event B occurring, given event A has occurred

* $P(A)$ - the Probability of event A

* $P(B)$ - the Probability of event B

One of the disadvantages of Naive - Bayes is that if you have no occurrences of a class label and a certain attribute value together, then the frequency - based Probability estimate will be zero. An approach to overcome this 'zero - frequency problem' in a Bayesian environment is to add one to the count for every attribute value - class combination when an attribute value doesn't occur with every class value.

Q5 → Section B

Backpropagation is a neural network learning algorithm. started by psychologist and neurobiologists to develop and test computational analogues of neurons.

A neural network is a set of connected input/output units where each connection has a weight associated with it. during the learning phase, the network learns by adjusting the weights so as to be able to predict the correct class label of the input tuples.

The weakness and the strength of neural network as a classifier are given below:-

Weakness

- long training time
- Require a number of parameters typically best determined empirically e.g. - the network topology or structure

- Poor interpretability : difficult to interpret the symbolic meaning behind the learned weights and of "hidden units" in the network.

Strength

- High tolerance to noisy data.

- Ability to classify untrained patterns

- Well-suited for continuous-valued input and outputs

- Successfully on a wide array of real world data

- Algorithms are inherently parallel

- Techniques have recently been developed for the extraction of rules from trained neural networks.

Q7 → Section C

* age

Age	Buy Computer	
	yes	NO
<=30	3 2	2 3
31-40	4	0
>40	3	2

* income

income	Buy Computer	
	yes	NO
high	2	2
medium	4	2
low	3	1

* Student

Student	Buy Computer	
	yes	NO
yes	6	1
no	3	4

* Credit rating

Credit rating	Buy Computer	
	yes	NO
fair	6	2
excellent	3	3

• $P(C_i)$:

$$P(\text{buys - computer} = \text{"yes"}) = 9/14 = 0.643$$

$$P(\text{buys - computer} = \text{"No"}) = 5/14 = 0.357$$

Compute $P(x|c_i)$ for each class

$$P(\text{age} = "<=30" | \text{buys-computer} = \text{"yes"}) = 2/9 = 0.222$$

$$P(\text{age} = "<=30" | \text{buys-computer} = \text{"NO"}) = 3/5 = 0.6$$

$$P(\text{income} = \text{"medium"} | \text{buys-computer} = \text{"yes"}) = 4/9 = 0.444$$

$$P(\text{income} = \text{"medium"} | \text{buys-computer} = \text{"NO"}) = 2/5 = 0.4$$

$$P(\text{student} = \text{"yes"} | \text{buys-computer} = \text{"yes"}) = 6/9 = 0.667$$

$$P(\text{student} = \text{"yes"} | \text{buys-computer} = \text{"NO"}) = 1/5 = 0.2$$

$$P(\text{credit-rating} = \text{"fair"} | \text{buys-computer} = \text{"yes"}) = 6/9 = 0.667$$

$$P(\text{credit-rating} = \text{"fair"} | \text{buys-computer} = \text{"NO"}) = 2/5 = 0.4$$

$x = (\text{age} <= 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$

$$\# P(x|c_i): P(x | \text{buys-computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$$

$$P(x | \text{buys-computer} = \text{"NO"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.19$$

$$\# P(x|c_i) \times P(c_i): P(x | \text{buys-computer} = \text{"yes"}) \times P(\text{buys-computer} = \text{"yes"}) = 0.044 \times 0.222 = 0.028$$

$$P(x | \text{buys-computer} = \text{"NO"}) \times P(\text{buys-computer} = \text{"NO"}) = 0.357 \times 0.19 = 0.007$$

Q8 → Section C

We know,

$$I_j = \sum_i w_{ij} O_i + \theta_j,$$

$$O_j = \frac{1}{1 + e^{-I_j}}$$

Net Input and Output Calculation:

Unit, j	Net Input I_j	Output O_j
4	$0.2 + 0 - 0.5 - 0.4 = -0.7$	$1/(1 + e^{0.7}) = 0.332$
5	$-0.3 + 0 + 0.2 + 0.2 = 0.1$	$1/(1 + e^{-0.1}) = 0.525$
6	$(-0.3)(0.332) - (0.2)(0.525) + 0.1 = 0.105$	$1/(1 + e^{0.105}) = 0.4$

Error at Each Node :

$$Error_j = O_j(1 - O_j)(T_j - O_j)$$

calculation of the error at Each Node

Unit, j	Error _j
4	$(0.332)(1 - 0.332)(0.1311)(-0.3) = -0.0087$
5	$(0.525)(1 - 0.525)(0.1311)(-0.2) = -0.0065$
6	$(0.474)(1 - 0.474)(1 - 0.0474) = 0.1311$

Weight and Bias updating

$$\Delta W_{ij} = (1) Error_j O_j$$

$$W_{ij} = W_{ij} + \Delta W_{ij}$$

$$\Delta \theta_j = (1) Error_j$$

$$\theta_j = -\theta_j + \Delta \theta_j$$

Weight
or
Bias

• New value

w_{46}

$$-0.3 + (0.9)(0.1311)(0.332) = -0.261$$

w_{56}

$$-0.2 + (0.9)(0.1311)(0.525) = -0.138$$

w_{14}

$$0.2 + (0.9)(-0.0087)(1) = 0.192$$

w_{15}

$$-0.3 + (0.9)(-0.0065)(1) = -0.306$$

w_{24}

$$0.4 + (0.9)(-0.0087)(0) = 0.4$$

w_{25}

$$0.1 + (0.9)(-0.0065)(0) = 0.1$$

w_{34}

$$-0.5 + (0.9)(-0.0087)(1) = -0.508$$

w_{35}

$$0.2 + (0.9)(-0.0065)(1) = 0.199$$

θ_6

$$0.1 + (0.9)(0.1311) = 0.218$$

θ_5

$$0.2 + (0.9)(-0.0065) = 0.194$$

θ_4

$$-0.4 + (0.9)(-0.0087) = -0.408$$