

## **BRAINWARE UNIVERSITY**

## Term End Examination 2020 - 21

Programme – Master of Computer Applications
Course Name – Machine Learning
Course Code - MCA501B

Semester / Year - Semester V

Time allotted: 85 Minutes

Full Marks: 70

[The figure in the margin indicates full marks. Candidates are required to give their answers in their own words as far as practicable.]

## Group-A

(Multiple Choice Type Question) 1 x 70=70

- 1. (Answer any Seventy)
- (i) The SVM's are less effective when:

a) The data is linearly separable

b) The data is clean and ready to use

c) The data is noisy and contains

d) None of these

overlapping points

- (ii) Suppose you are using RBF kernel in SVM with high Gamma value. What does this signify?
  - a) The model would consider even far away b) The model would consider only the points from hyperplane for modeling points close to the hyperplane for modeling
  - c) The model would not be affected by distance of points from hyperplane for modeling
- d) None of these
- (iii) The cost parameter in the SVM means:
  - a) The number of cross-validations to be made
- b) The kernel to be used
- c) The tradeoff between misclassification and simplicity of the model
- d) None of these
- (iv) Suppose you are building a SVM model on data X. The data X can be error prone which means that you should not trust any specific data point too much. Now think that you want to build a SVM model which has quadratic kernel

function of polynomial degree 2 that uses Slack variable C as one of it's hyper
parameter. Based upon that give the answer for following question. What
would happen when you use very large value of C(C->infinity)? Note: For
small C was also classifying all data points correctly

- a) We can still classify data correctly for given setting of hyper parameter C
- c) Can't Say

- b) We can't classify data correctly for given setting of hyper parameter C
- d) None of these
- (v) Which of the following are real world applications of the SVM?
  - a) Text and Hypertext Categorization
- b) Image Classification
- c) Clustering of News Articles
- d) All of these
- (vi) Which of the following option would you more likely to consider iterating SVM next time?
  - a) You want to increase your data points
- b) You want to decrease your data points
- c) You will try to calculate more variables
- d) You will try to reduce the features
- (vii) Suppose you gave the correct answer in previous question. What do you think that is actually happening? 1. We are lowering the bias 2. We are lowering the variance 3. We are increasing the bias 4. We are increasing the variance
  - a) 1 and 2

b) 2 and 3

c) 1 and 4

- d) 2 and 4
- (viii) In above question suppose you want to change one of it's(SVM) hyperparameter so that effect would be same as previous questions i.e model will not under fit?
  - a) We will increase the parameter C
- b) We will decrease the parameter C
- c) Changing in C don't affect
- d) None of these
- (ix) We usually use feature normalization before using the Gaussian kernel in SVM. What is true about feature normalization? 1. We do feature normalization so that new feature will dominate other 2. Sometimes, feature normalization is

helps when we use Gaussian kernel:	in SVM
a) 1	b) 1 and 2
c) 1 and 3	d) 2 and 3
	class classification problem and you want to nat you are using One-vs-all method. How M model in such case?
a) 1	b) 2
c) 3	d) 4
to train a SVM model on the data for Suppose you have same distribution	class classification problem and you want r that you are using One-vs-all method. of classes in the data. Now, say for training is taking 10 second. How many seconds nethod end to end?
a) 20	b) 40
c) 60	d) 80
to train a SVM model on the data for	class classification problem and you want r that you are using One-vs-all method. nes we need to train SVM in such case?
a) 1	b) 2
c) 3	d) 4
(xiii) What is/are true about kernel in dimensional data to high dimensional	-
a) 1	b) 2
c) 1 and 2	d) None of these
trees, individual weak learners are in	rue about boosting trees? 1. In boosting adependent of each other 2. It is the method ggregating the results of weak learners

not feasible in case of categorical variables 3. Feature normalization always

a) 1	b) 2
c) 3 and 4	d) 1 and 4
(xvi) In Random forest you can generate hund and then aggregate the results of these tree. W about individual(Tk) tree in Random Forest? I subset of the features 2. Individual tree is built tree is built on a subset of observations 4. Indiobservations	hich of the following is true  1. Individual tree is built on a t on all the features 3. Individual
a) 1 and 3	b) 1 and 4
c) 2 and 3	d) 2 and 4
(xvii) Which of the following algorithm doesn its hyperparameter? 1. Gradient Boosting 2. E. Random Forest	
a) 1 and 3	b) 1 and 4
c) 2 and 3	d) 2 and 4
(xviii) In random forest or gradient boosting a type. For example, it can be a continuous feature Which of the following option is true when yo features?	ure or a categorical feature.
a) Only Random forest algorithm handles real valued attributes by discretizing them	b) Only Gradient boosting algorithm handles real valued attributes by

b) 2

(xv) Which of the following is/are true about Random Forest and Gradient Boosting ensemble methods? 1. Both methods can be used for classification task 2. Random Forest is use for classification whereas Gradient Boosting is use for regression task 3. Random Forest is use for regression whereas Gradient

Boosting is use for Classification task 4. Both methods can be used for

d) None of these

a) 1

c) 1 and 2

regression task

1.		•	. 1
discret	117	7111A	them
aisci Ci	11/	ZIIIZ	uiciii

c) Both algorithms can handle real valued attributes by discretizing them	d) None of these
(xix) Which of the following is true about the Geach stage, introduce a new regression tree to contain model 2. We can use gradient decent manufaction	ompensate the shortcomings of
a) 1	b) 2
c) 1 and 2	d) None of these
(xx) The average squared difference between cl actual output.	assifier predicted output and
a) mean squared error	b) root mean squared error
c) mean absolute error	d) mean relative error
(xxi) A feed-forward neural network is said to b	be fully connected when
a) all nodes are connected to each other	b) all nodes at the same layer are connected to each other
c) all nodes at one layer are connected to all nodes in the next higher layer	d) all hidden layer nodes are connected to all output layer nodes
(xxii) How can you prevent a clustering algorith local optima?	nm from getting stuck in bad
a) Set the same seed value for each run	b) Use multiple random initializations
c) Both Set the same seed value for each runand Use multiple random initializations	d) None of these
(xxiii) To calculate the Median	
a) Middle value of samples	b) Arrange the samples in ascending order

d) All of these

c) Calculate middle position

(xxiv) The measures of dispersion include	
a) Range	b) Standard Deviation
c) Variance	d) None of these
(xxv) The range is	
a) Highest value-Lowest Value	b) Lowest Value- Highest value
c) Mean Value- Highest value	d) None of these
(xxvi) Variance is	
a) Sample mean of the squared deviations from the arithmetic mean	b) Arithmetic mean of the squared deviations from the sample mean
c) Sample mean of the squared deviations from the sample mean	d) None of these
(xxvii) Standard deviation is the	
a) Square of the variance	b) Cube of the variance
c) Square root of the variance	d) None of these
(xxviii) Coefficient of the correlation ranges fr	om
a) -1 to +1	b) 0 to +1
c) -1 to 0	d) None of these
(xxix) Chebysheff's theorem deals with	
a) Range	b) Variance
c) Standard deviation	d) None of these
(xxx) Conditional probability is related with	

(xxxi) Adding a non-important feature to a linear regression model may result

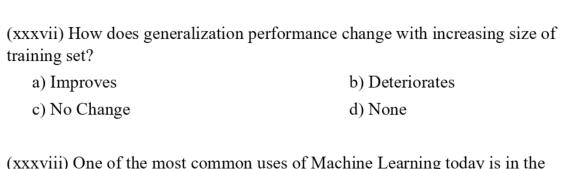
b) Chebysheff's Theorem

d) None of these

a) Naïve Bayes theorem

c) Pythagorian theorem

in. 1. Increase in R-square 2. Decrease	se in R-square
a) Only 1 is correct	b) Only 2 is correct
c) Either 1 or 2	d) None of these
1. Increase in K will result in higher Higher values of K will result in high	ns is/are true for K-fold cross-validation? time required to cross validate the result. 2. her confidence on the cross-validation K. 3. If K=N, then it is called Leave one number of observations.
a) 1 and 2	b) 2 and 3
c) 1 and 3	d) 1,2 and 3
(xxxiii) Which of the following option 5-fold cross validation with 10 differ	on is true for overall execution time for ent values of "max_depth"?
a) Less than 100 seconds	b) $100 - 300$ seconds
c) 300 – 600 seconds	d) More than or equal to 600 seconds
(xxxiv) What would you do in PCA	to get the same projection as SVD?
a) Transform data to zero mean	b) Transform data to zero median
c) Not possible	d) None of these
	ns can be used to get global minima in k- rithm for different centroid initialization 2. but the optimal number of clusters b) 1 and 3
	d) All of these
c) 1 and 2	
	yperparameters, higher value is better for of samples used for split 2. Depth of tree 3.
a) 1 and 2	b) 2 and 3
c) 1 and 3	d) Can't say



(xxxviii) One of the most common uses of Machine Learning today is in the domain of Robotics. Robotic tasks include a multitude of ML methods tailored towards navigation, robotic control and a number of other tasks. Robotic control includes controlling the actuators available to the robotic system. An example of this is control of a painting arm in automotive industries. The robotic arm must be able to paint every corner in the automotive parts while minimizing the quantity of paint wasted in the process. Which of the following learning paradigms would you select for training such a robotic arm?

a) Supervised learning

b) Unsupervised learning

c) Combination of supervised and unsupervised learning

d) Reinforcement learning

(xxxix) Let u be a  $n \times 1$  vector, such that u T u = 1. Let I be the  $n \times n$  identity matrix. The  $n \times n$  matrix A is given by (I ? kuuT), where k is a real constant. u itself is an eigenvector of A, with eigenvalue ?1. What is the value of k?

$$a) -2$$

b) -1

c) 2

d) 0

(xl) Consider the following 4 training examples • x = ?1, y = 0.0319 • x = 0, y = 0.8692 • x = 1, y = 1.9566 • x = 2, y = 3.0343 We want to learn a function f(x) = ax+b which is parametrized by (a, b). Using squared error as the loss function, which of the following parameters would you use to model this function.

a) 
$$(1, 1)$$

b) (1, 2)

d)(2, 2)

(xli) Which of the following are recommended applications of PCA?

a) To get more features to feed into a

b) Data compression: Reduce the

learning algorithm.

c) Preventing overfitting: Reduce the number of features (in a supervised learning linear regression: For most learning problem), so that there are fewer parameters to learn.

dimension of your data, so that it takes up less memory / disk space.

d) As a replacement for (or alternative to) applications, PCA and linear regression give substantially similar results.

(xlii) For which of the following problems would anomaly detection be a suitable algorithm?

- a) From a large set of primary care patient records, identify individuals who might have unusual health conditions.
- c) Given an image of a face, determine whether or not it is the face of a particular famous individual.
- b) Given data from credit card transactions, classify each transaction according to type of purchase (for example: food, transportation, clothing).
- d) From a large set of hospital patient records, predict which patients have a particular disease (say, the flu).

(xliii) Suppose you have trained an anomaly detection system for fraud detection, and your system that flags anomalies when p(x) is less than ?, and you find on the cross-validation set that it is missing many fraudulent transactions (i.e., failing to flag them as anomalies). What should you do?

a) Increase?

- b) Decrease?
- c) Both Increase? and Decrease?
- d) None of these

(xliv) What is the purpose of performing cross-validation?

- a) To assess the predictive performance of b) To judge how the trained model the models
  - performs outside the sample on test data
- c) Both To assess the predictive performance of the models and To judge how the trained model performs outside the sample on test data
- d) None of these

(xlv) Which of the following is a categorical outcome?

a) RMSE	b) RSquared
c) Accuracy	d) All of these
(xlvi) Which of the following statements about	regularization is not correct?
<ul><li>a) Using too large a value of lambda can cause your hypothesis to underfit the data.</li><li>c) Using a very large value of lambda cannot hurt the performance of your hypothesis.</li></ul>	<ul><li>b) Using too large a value of lambda can cause your hypothesis to overfit the data.</li><li>d) None of these</li></ul>
(xlvii) K-fold cross-validation is	
a) linear in K	b) quadratic in K
c) cubic in K	d) exponential in K
(xlviii) If I am using all features of my dataset my training set, but ~70% on validation set, when the control of the control	•
a) Underfitting	b) Nothing, the model is perfect
c) Overfitting	d) None of these
(xlix) Suppose you are using SVM with linear think that you increase the complexity (or degree What would you think will happen?	
<ul> <li>a) Increasing the complexity will overfit the data</li> </ul>	e b) Increasing the complexity will underfit the data
c) Nothing will happen since your model was already 100% accurate	d) None of these
(l) K-Means clustering algorithm is example of	Swhich model?
a) Connectivity models	b) Centroid models
c) Distribution models	d) Density Models
(li) Expectation-maximization algorithm is exa	mple of which model?

a) Connectivity models	b) Centroid models
c) Distribution models	d) Density Models
(E) DDCCAN and ODTICS are avanuals	of red ich and doll
(lii) DBSCAN and OPTICS are example	
a) Connectivity models	b) Centroid models
c) Distribution models	d) Density Models
(liii) In K Means Clustering algorithm, K	denotes
a) Number of associations	b) Number of regressions
c) Number of clusters	d) None of these
(liv) Hierarchical Clustering algorithm ter	minates when
a) there is only a single cluster left.	b) two nearest clusters are merged into the same cluster.
c) all the data points assigned to a clus their own	ster of d) None of these.
(lv) K means is	
a) Clustering algorithm	b) Classification algorithm.
c) Association algorithm	d) None of these.
(lvi) The time complexity of K Means is	
a) linear	b) quadratic
c) cubic	d) none of these
(lvii) The time complexity of hierarchical	clustering is
a) linear	b) quadratic
c) cubic	d) none of these
(lviii) Which metrics are used for deciding	g the closeness of two clusters?
a) Euclidean distance	b) Manhattan distance

c) Maximum distance	d) All of these	
(lix) Components of machine learning are		
a) Representation	b) Evaluation	
c) Optimization	d) All of these	
(lx) The sample of data used to fit the model is	called	
a) Training dataset	b) Validation dataset	
c) Test dataset	d) None of these	
(lxi) Suppose that a psychologist wants to evaluate the effectiveness of a new learning strategy. She randomly assigns students to two groups and assigns each student the same passage on a particular topic to study for half an hour. Subsequently each student participates in an individual assessment on the topic, where students of the one group use the new learning strategy, and students of the other group use any strategy they prefer. Which among the following is an extraneous variable in the above experiment?		
a) The choice of using two groups	b) The amount of time given to study the passage	
c) Existing knowledge about the passage among the students	d) The amount of time given to complete the assessment	
(lxii) I have trained a classifier, and to evaluate its performance I perform a 10-fold validation. I have obtained the following accuracies on the validation set in each of the runs - 0.90, 0.98, 0.95, 0.98, 0.97, 0.96, 0.94, 0.99, 0.96, 0.96. What is the sample standard deviation for the accuracies?		
a) 0.0243	b) 0.0256	
c) 0.000654439999999999	d) 0.000589	

(lxiii) A computer program is said to learn from experience E with respect to some task T and some performance measure P if its performance on T, as measured by P, improves with experience E. Suppose we feed a learning algorithm a lot of historical weather data, and have it learn to predict weather.

In this setting, what is E?

- a) The weather prediction task.
- b) The probability of it correctly predicting a future date's weather
- c) The process of the algorithm examining a large amount of historical weather data
- d) None of these.

(lxiv) Some of the problems below are best addressed using a supervised learning algorithm, and the others with an unsupervised learning algorithm. Which of the following would you apply supervised learning to? (Select all that apply.) In each case, assume some appropriate dataset is available for your algorithm to learn from.

- a) Given 50 articles written by male authors, and 50 articles written by female authors, learn to predict the gender of a new there are sub-types of spam mail manuscript's author (when the identity of this author is unknown).
  - b) Examine a large collection of emails that are known to be spam email, to discover if
- c) Given data on how 1000 medical patients d) All of these respond to an experimental drug (such as effectiveness of the treatment, side effects, etc.), discover whether there are different categories or "types" of patients in terms of how they respond to the drug, and if so what these categories are
- (lxv) Suppose that for some linear regression problem (say, predicting housing prices as in the lecture), we have some training set, and for our training set we managed to find some ?0, ?1 such that J(?0,?1)=0. Which of the statements below must then be true?
  - a) For this to be true, we must have ?0=0 and ?1=0 so that h?(x)=0

for every value of i=1,2,...,m.

- c) For this to be true, we must have y(i)=0
- b) Our training set can be fit perfectly by a straight line, i.e., all of our training examples lie perfectly on some straight line
  - d) We can perfectly predict the value of y even for new examples that we have not yet seen.

(lxvi) You run gradient descent for 15 iterations with ?=0.3 and compute J(?) after each iteration. You find that the value of J(?) decreases slowly and is still decreasing after 15 iterations. Based on this, which of the following conclusions seems most plausible?

- a) Rather than use the current value of ?, it'd be more promising to try a larger value of ? (say ?=1.0).
- it'd be more promising to try a smaller value of ? (say ?=0.1).

b) Rather than use the current value of?,

- c) ?=0.3 is an effective choice of learning rate.
- d) None of these

(lxvii) Which of the following are reasons for using feature scaling?

- a) It speeds up gradient descent by making it require fewer iterations to get to a good solution
- b) It is necessary to prevent gradient descent from getting stuck in local optima
- c) It speeds up gradient descent by making each iteration of gradient descent less expensive to compute.
- d) It prevents the matrix XTX (used in the normal equation) from being non-invertable (singular/degenerate)

(lxviii) Suppose Theta1 is a 5x3 matrix, and Theta2 is a 4x6 matrix. You set thetaVec = [Theta1(:), Theta2(:)]. Which of the following correctly recovers?

- a) reshape(thetaVec(16:39), 4, 6)
- b) reshape(thetaVec(15:38), 4, 6)
- c) reshape(thetaVec(16:24), 4, 6)
- d) reshape(thetaVec(15:39), 4, 6)

(lxix) Suppose you have a dataset with n = 10 features and m = 5000 examples. After training your logistic regression classifier with gradient descent, you find that it has underfit the training set and does not achieve the desired performance on the training or cross validation sets. Which of the following might be promising steps to take?

- a) Increase the regularization parameter?.
- b) Use an SVM with a Gaussian Kernel.
- c) Use an SVM with a linear kernel, without introducing new features.
- d) Reduce the number of examples in the training set

(lxx) Suppose you have implemented regularized logistic regression to classify

what object is in an image (i.e., to do object recognition). However, when you test your hypothesis on a new set of images, you find that it makes unacceptably large errors with its predictions on the new images. However, your hypothesis performs well (has low error) on the training set. Which of the following are promising steps to take?

- a) Try adding polynomial features.
- c) Try using a smaller set of features.
- b) Use fewer training examples.
- d) Try evaluating the hypothesis on a cross validation set rather than the test set.