

MACHINE LEARNING

In Q1 to Q5, only one option is correct, Choose the correct option:

1. In which of the following you can say that the model is overfitting?
 A) High R-squared value for train-set and High R-squared value for test-set.
 B) Low R-squared value for train-set and High R-squared value for test-set.
C) High R-squared value for train-set and Low R-squared value for test-set.
 D) None of the above
2. Which among the following is a disadvantage of decision trees?
 A) Decision trees are prone to outliers.
B) Decision trees are highly prone to overfitting.
 C) Decision trees are not easy to interpret
 D) None of the above.
3. Which of the following is an ensemble technique?
 A) SVM
C) Random Forest
 B) Logistic Regression
 D) Decision tree
4. Suppose you are building a classification model for detection of a fatal disease where detection of the disease is most important. In this case which of the following metrics you would focus on?
 A) Accuracy
C) Precision
 B) Sensitivity
 D) None of the above.
5. The value of AUC (Area under Curve) value for ROC curve of model A is 0.70 and of model B is 0.85. Which of these two models is doing better job in classification?
 A) Model A
B) Model B
 C) both are performing equal
 D) Data Insufficient

In Q6 to Q9, more than one options are correct, Choose all the correct options:

6. Which of the following are the regularization technique in Linear Regression??
 A) **Ridge**
 C) MSE
 B) R-squared
D) Lasso
7. Which of the following is not an example of boosting technique?
 A) Adaboost
B) Decision Tree
C) Random Forest
 D) Xgboost.
8. Which of the techniques are used for regularization of Decision Trees?
 A) Pruning
B) L2 regularization
C) Restricting the max depth of the tree
 D) All of the above
9. Which of the following statements is true regarding the Adaboost technique?
A) We initialize the probabilities of the distribution as $1/n$, where n is the number of data-points
B) A tree in the ensemble focuses more on the data points on which the previous tree was not performing well
 C) It is example of bagging technique
 D) None of the above

Q10 to Q15 are subjective answer type questions, Answer them briefly.

10. Explain how does the adjusted R-squared penalize the presence of unnecessary predictors in the model?

Answer: The adjusted R-squared is a modified version of R-squared that adjusts for predictors that are not significant in a regression model. Compared to a model with additional input variables, a lower adjusted R-squared indicates that the additional input variables are not adding value to the model.

MACHINE LEARNING

11. Differentiate between Ridge and Lasso Regression.

Answer:

	Lasso	Ridge
Type	L1 regularization	L2 regularization
Loss Function	Quadratic loss	Quadratic loss
Regularization mechanism	Model selection: by canceling some coefficients β_i ($ \beta_i = 0$)	Variable selection: by retaining all variables and constraining the parameter norm ($ \beta_i $)
Differentiability	Not differentiable	Differentiable (useful for gradient calculation)

12. What is VIF? What is the suitable value of a VIF for a feature to be included in a regression modelling?

Answer: The Variance Inflation Factor (VIF) is a measure of colinearity among predictor variables within a multiple regression. It is calculated by taking the ratio of the variance of all a given model's betas divide by the variance of a single beta if it were fit alone.

As a rule of thumb, a VIF of three or below is not a cause for concern. As VIF increases, the less reliable your regression results are going to be.

13. Why do we need to scale the data before feeding it to the train the model?

Answer: To ensure that the gradient descent moves smoothly towards the minima and that the steps for gradient descent are updated at the same rate for all the features, we scale the data before feeding it to the model.

14. What are the different metrics which are used to check the goodness of fit in linear regression?

Answer: MSE (Mean squared error), RMSE(Root mean squared error), or MAE(Mean Absolute error) are better be used to compare performance between different regression models.

15. From the following confusion matrix calculate sensitivity, specificity, precision, recall and accuracy.

Actual/Predicted	True	False
True	1000	50
False	250	1200

Answer:

- Sensitivity = $TP / (TP + FN) = 1000 / (1000 + 50) = 0.952$
- Specificity = $TN / (FP + TN) = 1200 / (250 + 1200) = 0.827$
- Precision = $TP / (TP + FP) = 1000 / (1000 + 250) = 0.8$
- Recall = $TP / (TP + FN) = 1000 / (1000 + 50) = 0.952$
- Accuracy = $(TP + TN) / (TP + TN + FP + FN) = (1000 + 1200) / (1000 + 1200 + 50 + 250) = 0.88$