

MACHINE LEARNING

(ASSIGNMENT - Worksheet6 Answers – Vivek Kumar Sahu – Internship 35)

(Marked answers in Bold)

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
C) Random Forest
B) **Decision Tree**
D) Xgboost.
8. Which of the techniques are used for regularization of Decision Trees?
A.) **Pruning**
C) Restricting the max depth of the tree
B) L2 regularization
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

MACHINE LEARNING

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?

Ans. The adjusted R-squared is a modified version of R-squared that has been adjusted for the number of predictors in the model. The adjusted R-squared increases only if the new term improves the model more than would be expected by chance. It decreases when a predictor improves the model by less than expected by chance. The adjusted R-squared can be negative, but it's usually not. It is always lower than the R-squared.

11. Differentiate between Ridge and Lasso Regression.

Ans.

Lasso Regression :

Lasso regression stands for Least Absolute Shrinkage and Selection Operator. It adds penalty term to the cost function. This term is the absolute sum of the coefficients. As the value of coefficients increases from 0 this term penalizes, cause model, to decrease the value of coefficients in order to reduce loss. The difference between ridge and lasso regression is that it tends to make coefficients to absolute zero as compared to Ridge which never sets the value of coefficient to absolute zero.

Ridge Regression :

In Ridge regression, we add a penalty term which is equal to the square of the coefficient. The L2 term is equal to the square of the magnitude of the coefficients. We also add a coefficient λ to control that penalty term. In this case if λ is zero then the equation is the basic OLS else if $\lambda > 0$ then it will add a constraint to the coefficient. As we increase the value of λ this constraint causes the value of the coefficient to tend towards zero. This leads to tradeoff of higher bias (dependencies on certain coefficients tend to be 0 and on certain coefficients tend to be very large, making the model less flexible) for lower variance.

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

Ans. The Variance Inflation Factor (VIF) measures the severity of multicollinearity in regression analysis. It is a statistical concept that indicates the increase in the variance of a regression coefficient as a result of collinearity.

It explains multicollinearity. In general, a VIF above 10 indicates high correlation and is cause for concern.

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

Ans 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. Having features on a similar scale can help the gradient descent converge more quickly towards the minima.

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

Ans: The sum of squares due to error (SSE)

- R-square
- Adjusted R-square
- Root mean squared error (RMSE)

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

Actual/Predicted	True	False
True	1000	50

MACHINE LEARNING

False	250	1200
-------	-----	------

Ans.

Sensitivity = $1000/1050 = 0.95$

Specificity = $1200/1450 = 0.83$

Precision = $1000/1250 = 0.8$

Recall = $1000/1050 = 0.95$

Accuracy = $2200/2500 = 0.88$

MACHINE LEARNING