# MACHINE LEARNING

1.  Which of the following in sk-learn library is used for hyper parameter tuning?
    A) GridSearchCV()                 B) RandomizedCV()
    C) K-fold Cross Validation        D) All of the above
    ANSWER: GridSearchCV()

2.  In which of the below ensemble techniques trees are trained in parallel?
    A) Random forest                  B) Adaboost
    C) Gradient Boosting              D) All of the above
    ANSWER: Random forest

3.  In machine learning, if in the below line of code:
    *sklearn.svm.SVC (C=1.0, kernel='rbf', degree=3)*
    we increasing the C hyper parameter, what will happen?
    A) The regularization will increase      B) The regularization will decrease
    C) No effect on regularization           D) kernel will be changed to linear
    ANSWER: The regularization will decrease

4.  Check the below line of code and answer the following questions:
    *sklearn.tree.DecisionTreeClassifier(\*criterion='gini',splitter='best',max_depth=None, min_samples_split=2)*
    Which of the following is true regarding max_depth hyper parameter?
    A) It regularizes the decision tree by limiting the maximum depth up to which a tree can be grown.
    B) It denotes the number of children a node can have.
    C) both A & B
    D) None of the above
    ANSWER: both A & B

5.  Which of the following is true regarding Random Forests?
    A) It's an ensemble of weak learners.
    B) The component trees are trained in series
    C) In case of classification problem, the prediction is made by taking mode of the class labels predicted by the component trees.
    D) None of the above
    ANSWER: It's an ensemble of weak learners

6.  What can be the disadvantage if the learning rate is very high in gradient descent?
    A) Gradient Descent algorithm can diverge from the optimal solution.
    B) Gradient Descent algorithm can keep oscillating around the optimal solution and may not settle.
    C) Both of them
    D) None of them
    ANSWER: Gradient Descent algorithm can diverge from the optimal solution.

7.  As the model complexity increases, what will happen?
    A) Bias will increase, Variance decrease      B) Bias will decrease, Variance increase
    C)both bias and variance increase             D) Both bias and variance decrease.
    ANSWER: Bias will decrease, Variance increase

# MACHINE LEARNING

8. Suppose I have a linear regression model which is performing as follows:
   Train accuracy=0.95 and Test accuracy=0.75
   Which of the following is true regarding the model?
   A) model is underfitting          B) model is overfitting
   C) model is performing good      D) None of the above
   ANSWER: model is overfitting

**Q9 to Q15 are subjective answer type questions, Answer them briefly.**

9. Suppose we have a dataset which have two classes A and B. The percentage of class A is 40% and percentage of class B is 60%. Calculate the Gini index and entropy of the dataset.
   ANSWER:

10. What are the advantages of Random Forests over Decision Tree?
    ANSWER: avoids and prevents overfitting by using multiple trees

11. What is the need of scaling all numerical features in a dataset? Name any two techniques used for scaling.

    ANSWER:To get better prediction data need to be scaled and it can be done by standard scaling or minmax scaling.

12. Write down some advantages which scaling provides in optimization using gradient descent algorithm.

ANSWER: Gradient descent is an optimization algorithm used to find the values of parameters (coefficients) of a function (f) that minimizes a cost function (cost).

Gradient descent is best used when the parameters cannot be calculated analytically (e.g. using linear algebra) and must be searched for by an optimization algorithm.

13. In case of a highly imbalanced dataset for a classification problem, is accuracy a good metric to measure the performance of the model. If not, why?
    ANSWER: The most common metric used to evaluate the performance of a classification predictive model is classification accuracy. Typically, the accuracy of a predictive model is good (above 90% accuracy), therefore it is also very common to summarize the performance of a model in terms of the error rate of the model.

14. What is "f-score" metric? Write its mathematical formula.
    ANSWER: An F-score is the harmonic mean of a system's precision and recall values. It can be calculated by the following formula: **2 x [(Precision x Recall) / (Precision + Recall)]**

15. What is the difference between fit(), transform() and fit_transform()?
    ANSWER: The fit(data) method is used to compute the mean and std dev for a given feature to be used further for scaling. The transform(data) method is used to perform scaling using mean and std dev calculated using the . fit() method. The fit_transform() method does both fits and transform