In [ ]:

```
1  #*******************************************************************************
2  #                      Machine Learning 1 : Assignement
3  #
4  #*******************************************************************************
5  # 1. What are the three stages to build the hypotheses or model in machine learning?
6  # 2.   What is the standard approach to supervised learning?
7  # 3.   What is Training set and Test set?
8  # 4.   What is the general principle of an ensemble method and what is bagging and boosting in ensemble method?
9  # 5.   How can you avoid overfitting ?
10
11
```

```
In [ ]:    1  1. What are the three stages to build the hypotheses or model in machine learning?
           2
           3  The Three Stages for building the hypothesis are as follows
           4
           5  a. Model Building /DataEngineering : These include following
           6      1. Understanding the Business /Domain ,This will help to understand and identify correctly the
           7         problem that is occuring ,that needs to be solved and build a case for the same.
           8      2. Once the problem is identified , then identify the data sources from which data can be obtained
           9         for solving the problem, the Data can be structured data as databases or unstructured like emails,
          10         documents,pictures etc.
          11
          12         Once the Data is obtained , then preprocess/normalize the Data by processes like Data Cleaning
          13         (filling missing values, ...), Data integration, data transformation.
          14         The normalized data is then separated into Training Set and Testing Set.
          15
          16      3. Developing or selecting a model/Algorithim , which helps to solve the use case .For this The
          17         training Data Set is used and model is build is build for it .
          18         Once the model is built , it is validated with the Validate set.
          19
          20
          21  b. Model Testing/Data Intelligence :
          22      In these stage,You test the buildModel with the testing Data Set, several times  and analyse the predicted
          23      results , if they are within the predicted score then model is performing as needed , but if
          24      accuracy of predicted score is not good , then identify the problem ,rectify,finetune the model to
          25      improve score , else some other algorithms are tried .
          26
          27  c. Model Appling/Deployment :
          28      In this stage , once the model is working as predicted and suffcient accuracy, then apply or deploy
          29      the model in the business environment  .
          30
          31
          32
          33  2. What is the standard approach to supervised learning?
          34
          35      The Standard approach for Supervised learning is :
          36      PreProcess the Raw Data and transform it into Data Features .The Data is split into training and
          37      Test Data . The Algorithm is then trained with the training data and labels for these DataFeatures.
          38      The Model is tuned to get the better accuracy for different algoritm parameters.
          39      The Fine tuned Model is then tested to predict Labels for the Test Data.
          40      Based on the Training  it has achieved , it appropriately labels the new Data.
          41
```

```
42   3. What is Training set and Test set?
43
44      The Raw Data is preprocessed by Data Cleaning, Data integretion , Data transformation processes like normalizatio
45      The Processed Data is divided into Two sets , Training and Test Set in the ratio 70:30 percent.
46
47      Training Set is the Data that is used for training the model/algorithm and hence the name.
48      The training set includes labels that by which the Model learns their mapping and this will help
49      it to make predictions of labels  for the new Data.
50
51      Once the Learning by the model is done , the Test set (the subset of input normalised data) is used
52      to test the model and predict the labels for the new data. The predictions made by the test Set
53      helps to give the accuracy of the predictions done by the Model.
54
55
56   4. What is the general principle of an ensemble method and what is bagging and boosting in ensemble method?
57
58      The general principle of ensemble is , it is supervised meta machine learning model which combines
59      the predictions of many small models to improve on the accuracy of predictions and make a robust ,
60      generalized model.
61
62      Ensemble is achieved by methods like bagging , boosting , stacking , to address the issues like statistical, comp
63      and representational respectively
64
65      Bagging is a type of ensemble where many models/estimators are build independently on subsets of the data
66      and then their predictions are averaged. The result /prediction of bagging model is better /more accurate
67      then each of the sub models. This method addresses the statistical issue faced by ensemble model by limiting
68      the variance between the data points without much affecting the bias.
69      eg is random Forest
70
71      Boosting is a type of ensemble where model/estimators are build sequentially and thesubsequent  model
72      works on the weaknesses of the previous model to make a better predicting model i.e by combining many
73      weak models it makes a stronger model .This method addresses the issue of computational faced by reducing
74      the bias of the data points without affecting the variance between the data points.
75      eg. Adaboost
76
77   5. How can you avoid overfitting ?
78
79      Overfitting in the model occurs when we train it with lot of data , in doing so , the model not only learns from
80      but also from the noise or inaccuracies in the data . It starts considering the errors as the actual data and
81      hence becomes overfitted , resulting in wrongly predicting new data .
82      To avoid overfitting , one needs to study the behaviours of data , if it is  linear , then use linear algoritm
83      and when non linear , then use paramters in the model classifier like max_features, max_dept, min_samples_leaf
```

```
84    to limit overfitting
85
```

In [ ]:    1