

Q What is Ensemble Trees.

Ans:

It combines a number of decision trees in order to reduce the risk of overfitting.

Q Write Ensemble Tree Algorithm

Ans:

Input: Training Data, D , C_{4.5} Learning Algorithm

Output: A set of trees, DT^*

Method:

1. $DT^* = \phi$

2. Create sub-datasets, $D_1, \dots, D_i, \dots, D_n$ from the training data, D ;

3. for $i=1$ to k do

 group features in D_i into m groups;

4.

 for $j=1$ to m do

5.

6. build a DT_j with j th feature group;

7.

 compute $error(DT_j)$ on D_i ;

8.

 if $error(DT_j) \leq \text{threshold-value}$ then

9.

$DT^* = DT^* \cup DT_j$;

10.

 end if

11.

 end for

12 end for

Q Define Class Imbalanced Problem

Page-761

Ans:

It is the Problem in machine learning where the total number of a class of data (positive) is far less than total numbers of another class of data (negative).

Q Write the Causes of overfitting & Overlapping of classes

Ans:

The causes of overfitting & overlapping are

- a. Noise
- b. High dimension
- c. Small Sample size.

Q Which methods are used for class imbalancing

Ans:

Three methods are used in class imbalancing. They are

- a. Sampling Method :
 - i. Under Sampling
 - ii. Over Sampling
- b. Cost - Sensitive Learning Method
- c. Ensemble methods :
 - i. Bagging
 - ii. Boosting

Q Describe Under-Sampling & Over Sampling

Undersampling :

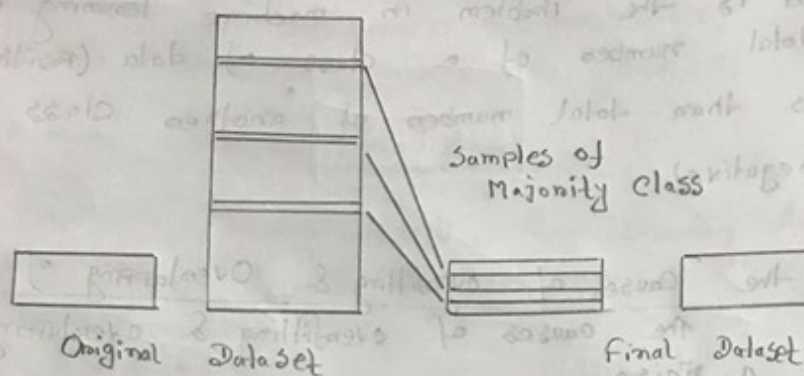


Figure: Undersampling

By Undersampling

- Run-time can be improved by decreasing the amount of training dataset.
- Helps to solve the memory problem.

It causes information loss.

Over Sampling :

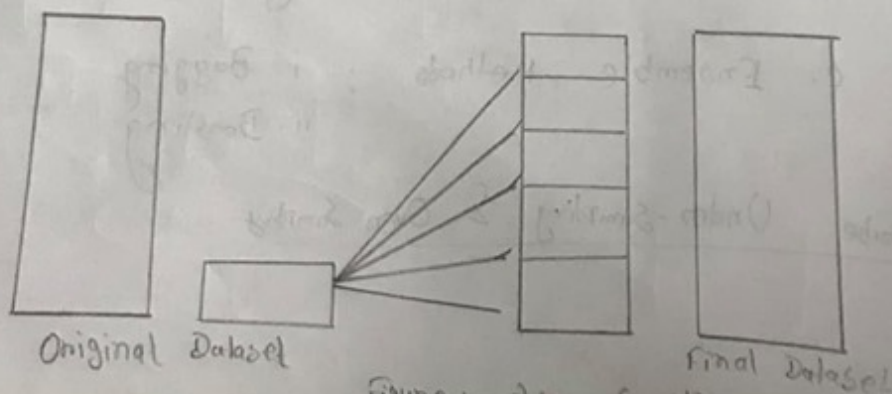


Figure: Over Sampling

By Oversampling

- a. No loss of information
- b. Mitigate overfitting caused by oversampling

Q What is Cost Sensitive Learning. Write the advantage

Ans:

It takes the mis-classification costs into consideration by minimising the total cost. The goal of this technique is to ~~Pursue~~ gain a high accuracy of classifying examples into a set of known classes.

Advantage

- a. Avoids Pre-selection of Parameters.
- b. Auto adjust decision hyperline.

Explain Data Balancing Method

Ans:

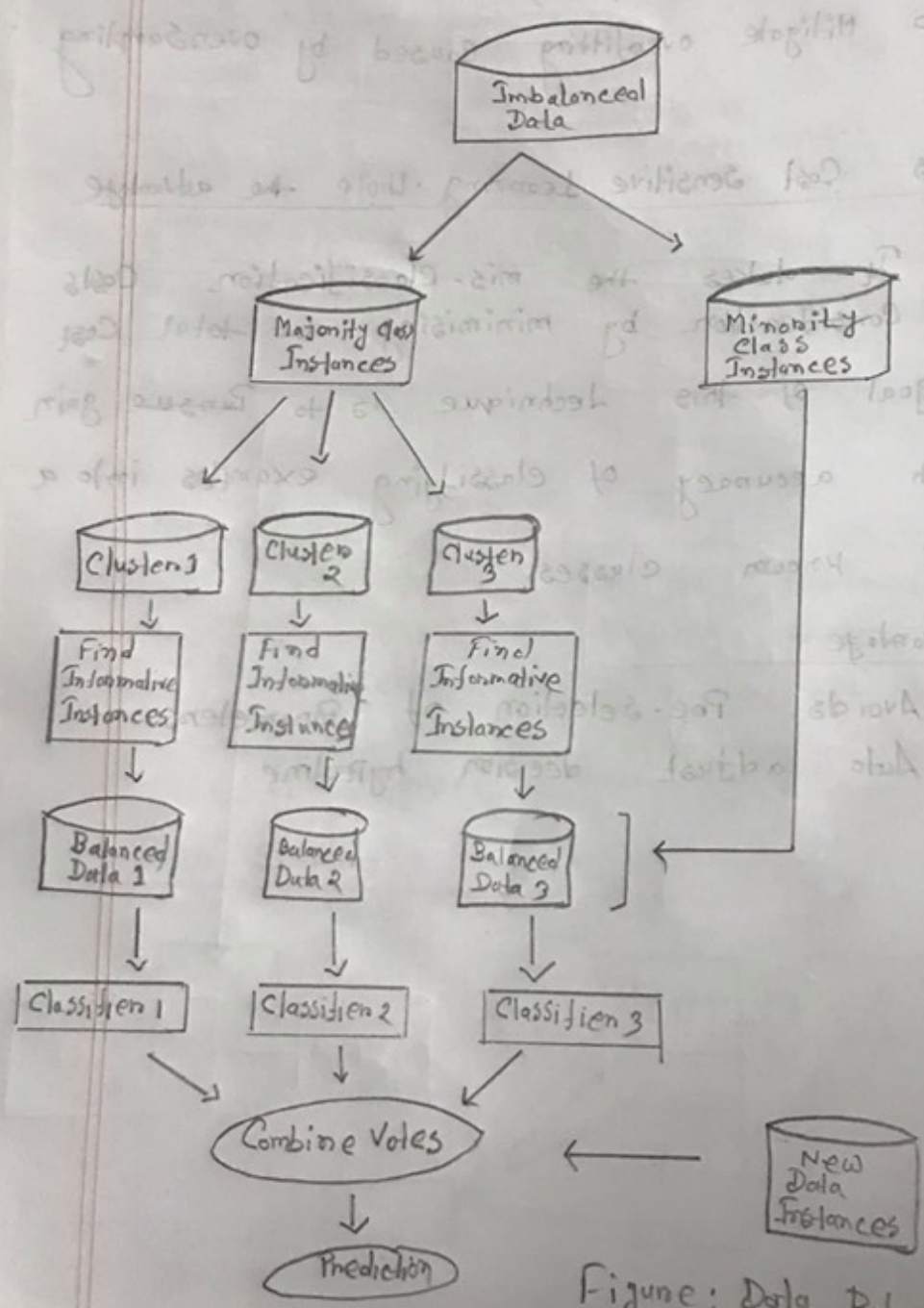


Figure: Data Balancing

For Data Balancing

- a. Majority class instances are clustered into several clusters.
- b. Next comes finding most informative instances in each cluster.
- c. Informative instances are close to the center or border of the cluster.
- d. With these most informative instances several datasets are created.
- e. Every dataset should have almost equal numbers of minority - majority classes instances.
- f. Multiple classifiers are trained using these datasets

Q. Explain Active Learning.

Ans.

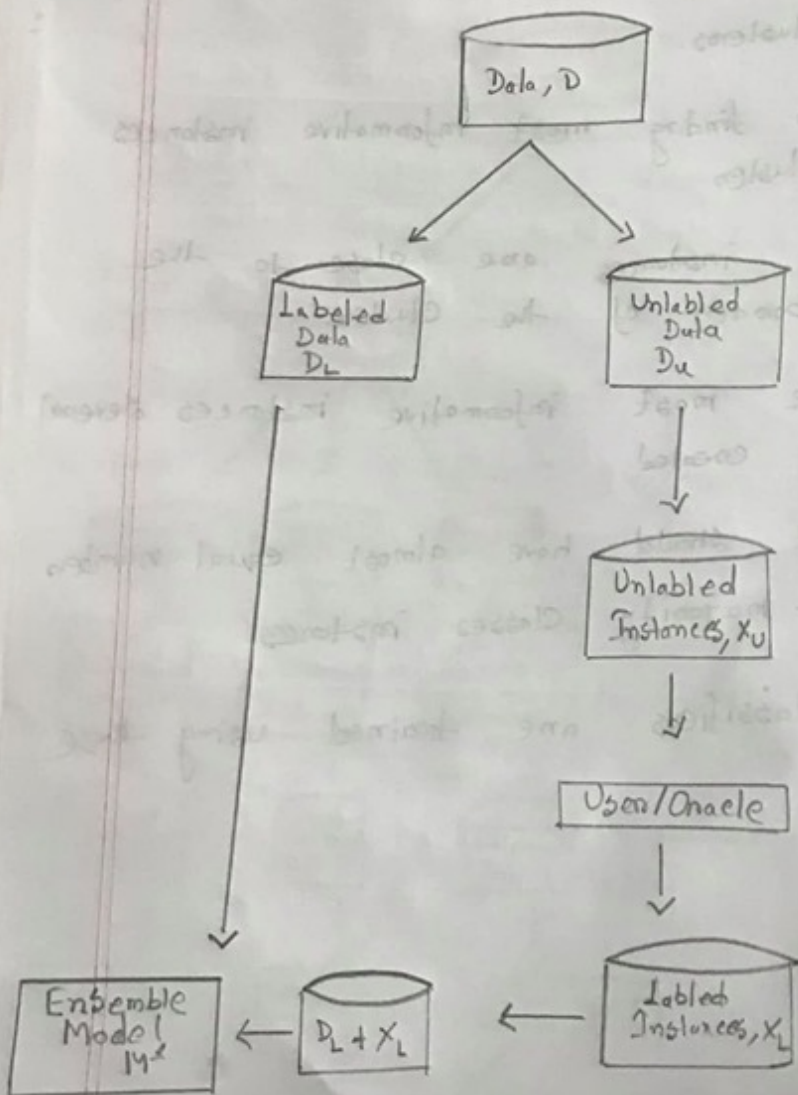


Figure: Active Learning

An active learning can ask a user to label an instance, which can be from a set of unlabelled instances.

Naive Bayes classifier & clustering are used to find the most informative instances for labeling as part of active learning. It follows two key rules

- a. Instances close to centers of clusters & borders of clusters.
- b. Posterior probabilities of instances are equal / very close.

Q. Write the steps of Data Pre-Processing

Ans:

The steps of Data Pre-Processing are:

a. Data Cleaning: Dealing with missing values.

b. Data Integration: Merge multiple sources data into coherent data store.

c. Data Transformation: It includes four

sub steps

i. Normalisation

ii. Aggregation

iii. Generalisation

iv. Feature Construction.

d. Data Reduction:

features.

Eliminating redundant

e. Data Discretisation:

Reduction of a number of values of a continuous feature.

Q What is feature selection. Write three reasons of feature selection

Ans:

Feature selection means the process of selecting a subset of relevant features

The three reasons of feature selection are

- Simplification of models
- Shorten training times
- Reducing overfitting.

Q What is Entropy

Ans:

The measure of randomness in the information being processed.

The End

30/08/2019