



COURSE LABORATORY MANUAL

v_{MAP} , given the attribute values (a_1, a_2, \dots, a_n) that describe the instance.

$$v_{MAP} = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2, \dots, a_n)$$

We can use Bayes theorem to rewrite this expression as

$$\begin{aligned} v_{MAP} &= \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)} \\ &= \operatorname{argmax}_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j) P(v_j) \end{aligned}$$

Now we could attempt to estimate the two terms in Equation (19) based on the training data. It is easy to estimate each of the $P(v_j)$ simply by counting the frequency with which each target value v_j occurs in the training data.

The naive Bayes classifier is based on the simplifying assumption that the attribute values are conditionally independent given the target value. In other words, the assumption is that given the target value of the instance, the probability of observing the conjunction a_1, a_2, \dots, a_n , is just the product of the probabilities for the individual attributes: $P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$.

Substituting this, we have the approach used by the naive Bayes classifier.

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

where v_{NB} denotes the target value output by the naive Bayes classifier.

When dealing with continuous data, a typical assumption is that the continuous values associated with each class are distributed according to a Gaussian distribution. For example, suppose the training data contains a continuous attribute, x . We first segment the data by the class, and then compute the mean and variance of x in each class.

Let μ be the mean of the values in x associated with class C_k , and let σ_k^2 be the variance of the values in x associated with class C_k . Suppose we have collected some observation value v . Then, the probability distribution of v given a class C_k , $p(x=v | C_k)$ can be computed by plugging v into the equation for a Normal distribution parameterized by μ and σ_k^2 . That is

$$p(x = v | C_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(v-\mu_k)^2}{2\sigma_k^2}}$$

Above method is adopted in our implementation of the program.

Pima Indian diabetes dataset

This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset.



COURSE LABORATORY MANUAL

6. PROCEDURE / PROGRAMME :

```
import csv, random, math
import statistics as st

def loadCsv(filename):
    lines = csv.reader(open(filename, "r"));
    dataset = list(lines)
    for i in range(len(dataset)):
        dataset[i] = [float(x) for x in dataset[i]]
    return dataset

def splitDataset(dataset, splitRatio):
    testSize = int(len(dataset) * splitRatio);
    trainSet = list(dataset);
    testSet = []
    while len(testSet) < testSize:
        #randomly pick an instance from training data
        index = random.randrange(len(trainSet));
        testSet.append(trainSet.pop(index))
    return [trainSet, testSet]

#Create a dictionary of classes 1 and 0 where the values are the
#instancs belonging to each class

def separateByClass(dataset):
    separated = {}
    for i in range(len(dataset)):
        x = dataset[i]
        if (x[-1] not in separated):
            separated[x[-1]] = []
        separated[x[-1]].append(x)
    return separated

def compute_mean_std(dataset):
    mean_std = [ (st.mean(attribute), st.stdev(attribute))
                 for attribute in zip(*dataset)]; #zip(*res) transposes a matrix (2-d array/list)
    del mean_std[-1] # Exclude label
    return mean_std

def summarizeByClass(dataset):
    separated = separateByClass(dataset);
    summary = {} # to store mean and std of +ve and -ve instances
    for classValue, instances in separated.items():
        #summaries is a dictionary of tuples(mean,std) for each class value
        summary[classValue] = compute_mean_std(instances)
    return summary

#For continuous attributes p is estimated using Gaussian distribution
def estimateProbability(x, mean, stdev):
    exponent = math.exp(-(math.pow(x-mean,2)/(2*math.pow(stdev,2))))
    return (1 / (math.sqrt(2*math.pi) * stdev)) * exponent
```



COURSE LABORATORY MANUAL

```
def calculateClassProbabilities(summaries, testVector):
    p = {}
    #class and attribute information as mean and sd
    for classValue, classSummaries in summaries.items():
        p[classValue] = 1
        for i in range(len(classSummaries)):
            mean, stdev = classSummaries[i]
            x = testVector[i] #testvector's first attribute
            #use normal distribution
            p[classValue] *= estimateProbability(x, mean, stdev);
    return p

def predict(summaries, testVector):
    all_p = calculateClassProbabilities(summaries, testVector)
    bestLabel, bestProb = None, -1
    for lbl, p in all_p.items():#assigns that class which has he highest prob
        if bestLabel is None or p > bestProb:
            bestProb = p
            bestLabel = lbl
    return bestLabel

def perform_classification(summaries, testSet):
    predictions = []
    for i in range(len(testSet)):
        result = predict(summaries, testSet[i])
        predictions.append(result)
    return predictions

def getAccuracy(testSet, predictions):
    correct = 0
    for i in range(len(testSet)):
        if testSet[i][-1] == predictions[i]:
            correct += 1
    return (correct/float(len(testSet))) * 100.0

dataset = loadCsv('data51.csv');
print('Pima Indian Diabetes Dataset loaded...')
print('Total instances available :',len(dataset))
print('Total attributes present :',len(dataset[0])-1)

print("First Five instances of dataset:")
for i in range(5):
    print(i+1 , ':' , dataset[i])

splitRatio = 0.2
trainingSet, testSet = splitDataset(dataset, splitRatio)
print('\nDataset is split into training and testing set.')
print('Training examples = {0} \nTesting examples = {1}'.format(len(trainingSet),
                                                                len(testSet)))

summaries = summarizeByClass(trainingSet);
predictions = perform_classification(summaries, testSet)

accuracy = getAccuracy(testSet, predictions)
print('\nAccuracy of the Naive Bayesian Classifier is :', accuracy)
```



COURSE LABORATORY MANUAL

7. RESULTS & CONCLUSIONS:

Sample Result

Pima Indian Diabetes Dataset loaded...

Total instances available : 768

Total attributes present : 8

First Five instances of dataset:

1 : [6.0, 148.0, 72.0, 35.0, 0.0, 33.6, 0.627, 50.0, 1.0]

2 : [1.0, 85.0, 66.0, 29.0, 0.0, 26.6, 0.351, 31.0, 0.0]

3 : [8.0, 183.0, 64.0, 0.0, 0.0, 23.3, 0.672, 32.0, 1.0]

4 : [1.0, 89.0, 66.0, 23.0, 94.0, 28.1, 0.167, 21.0, 0.0]

5 : [0.0, 137.0, 40.0, 35.0, 168.0, 43.1, 2.288, 33.0, 1.0]

Dataset is split into training and testing set.

Training examples = 615

Testing examples = 153

Accuracy of the Naive Bayesian Classifier is : 73.85

8. LEARNING OUTCOMES :

- The student will be able to apply naive bayesian classifier for the relevent problem and analyse the results.

9. APPLICATION AREAS:

- Real time Prediction: Naive Bayes is an eager learning classifier and it is sure fast. Thus, it could be used for making predictions in real time.
- Multi class Prediction: This algorithm is also well known for multi class prediction feature. Here we can predict the probability of multiple classes of target variable.
- Text classification/ Spam Filtering/ Sentiment Analysis: Naive Bayes classifiers mostly used in text classification (due to better result in multi class problems and independence rule) have higher success rate as compared to other algorithms. As a result, it is widely used in Spam filtering (identify spam e-mail) and Sentiment Analysis (in social media analysis, to identify positive and negative customer sentiments)
- Recommendation System: Naive Bayes Classifier and Collaborative Filtering together builds a Recommendation System that uses machine learning and data mining techniques to filter unseen information and predict whether a user would like a given resource or not

10. REMARKS: