Report: Naive Bayes Classifier

Solution 1: Naïve Bayes Classifier for UCI Census-Income (KDD) Data Set

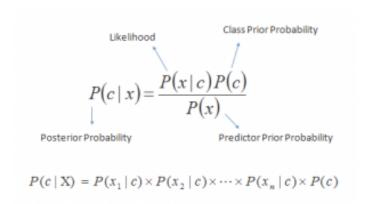
Pre Processing of the data:

- The data contains 37 nominal features, used as categorical features for classification.
- There are 7 continuous features, out of which age, wage per hour, capital gains were binned according to the histograms plotted.
- The missing data for categorical features was replaced with the mode of the feature as mean cannot be used and it was seen that one of the resulting values was dominant in the observations for the features.

```
== " main ":
train data = read data("dataset/census-income.data")
test data = read data("dataset/census-income.data")
print(train data.apply(num missing, axis=0)) #axis=0 defines that function is to be applied on each column
#plt.hist(list(train_data['age']), facecolor='green', alpha=0.75, bins=9)
bins = [0, 10, 20, 30, 40, 50, 60, 70, 80, 90]
group_names = [1,2,3,4,5,6,7,8,9]
categories = pd.cut(train data['age'], bins, labels=group names)
train data['age'] = pd.cut(train data['age'], bins, labels=group names)
categories = pd.cut(test_data['age'], bins, labels=group_names)
test data['age'] = pd.cut(test data['age'], bins, labels=group names)
train_data.hist(column="wage_per_hour",bins=30)
group_names = [1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30]
train_data['wage_per_hour'] = pd.qcut(train_data['wage_per_hour'], 30, labels=group_names)
train_data.hist(column="capital gains",bins=30)
group names = [1,2,3,4,5,6,7,8,\bar{9},10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30]
train data['capital gains'] = pd.qcut(train data['capital gains'], 30, labels=group names)
cat columns = train data.select dtypes(['object']).columns
train_data[cat_columns] = train_data[cat_columns].apply(lambda x: x.astype('category'))
train_data[cat_columns] = train_data[cat_columns].apply(lambda x: x.cat.codes)
test_data[cat_columns] = test_data[cat_columns].apply(lambda x: x.astype('category'))
test_data[cat_columns] = test_data[cat_columns].apply(lambda x: x.cat.codes)
train data[cat columns] = coding(train data[cat columns], {'N':0,'Y':1})
train data.to csv("dataset/train data.csv", sep=',')
test_data.to_csv("dataset/test_data.csv", sep=',')
```

Naive Bayes Classifier:

The Naive Bayes Classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. It is easy to build and useful for large datasets.



Pipeline of the classifier:

```
_ == " main ":
name
accuracy_mat = []
for t in range(10):
    model = Model("dataset/train_data.csv")
   model.load_database()
    splitRatio = 0.7
    model.k_fold_split(splitRatio)
   print('Size of train=',len(model.trainingSet),' and test=',len(model.testSet))
   model.summarize by class()
    predictions = model.test classifier()
    accuracy = model.get_accuracy(predictions)
    print('Accuracy[',t,']: ',accuracy)
    accuracy_mat.append(accuracy)
mean_accuracy = np.mean(accuracy_mat)
stddev_accuracy = np.std(accuracy_mat, ddof = 1)
print('For 10 runs, Mean Accuracy : ',mean_accuracy,' and Standard Deviation: ',stddev_accuracy)
```

Test Runs:

Size of train= 139666 and test= 59857

• The data was sampled as 70% train and 30% test for 10 runs and the following accuracy and standard deviation was reported:

```
Accuracy[ 0 ]: 70.62164826169037
Size of train= 139666
                      and test= 59857
Accuracy[ 1 ]:
              70.78370115441803
Size of train= 139666
                      and test= 59857
Accuracy[ 2 ]: 70.71019262575805
Size of train= 139666
                      and test= 59857
Accuracy[ 3 ]: 70.60995372303991
Size of train= 139666
                      and test= 59857
Accuracy[ 4 ]: 70.5297626008654
Size of train= 139666
                      and test= 59857
Accuracy[ 5 ]: 70.79205439631122
Size of train= 139666
                      and test= 59857
Accuracy[ 6 ]: 70.93907145363116
Size of train= 139666
                      and test= 59857
Accuracy[ 7 ]: 70.62666020682626
Size of train= 139666
                      and test= 59857
Accuracy[ 8 ]: 70.63668409709808
Size of train= 139666 and test= 59857
Accuracy[ 9 ]: 70.66174382277762
For 10 runs, Mean Accuracy: 70.6911472342 and Standard Deviation: 0.118337790033
```

Solution 2: Derivations

we want to find P(Y X) where Y= label, X= data
but dince P(XIY) is easier to model, we write
but since $P(X Y)$ is easier to model, we write $P(Y X) = P(X Y)P(Y)$
POSterior PLX)
pix) is unknown but has a known parametric form,
P(X/O) Sunknown parameter.
Sunknown paramew.
we night have empert domain knowledge of set of prior > P(0) & posterior > P(0/D) samples.
.: x, D selected independently => P(x/0, D): P(x/0)
=> P(x/D) = Sp(x/0)p(0/D)do
distribution:
for gaussian distribution: $p(x D), p(\theta D) = ?$ where $p(x \mu) \sim N(\mu, \Sigma)$
La constitución de la constituci
p(x/M) ~ N(µ, 62) known assuming p(µ) ~ N(µ0, 52)
p(x) (x) (x) best suess for a uncertainty
2 .12 0(D) 38(D) - XE(D) D) XIV
then Bayes formula => P(µ(D) = P(D) P(µ) = XP(D) P(µ)
Spippu) phidm (4)
hunce we have, (D, (3, (4) =>)
P(UID) = QX [1 e -1/2 (x-11)2] [1 e -6)

$$\begin{array}{c} - \frac{1}{2} \left(\frac{2(\mu - \chi_{1})^{2}}{\sigma^{2}} + \frac{(\mu - \mu_{1})^{2}}{\sigma^{2}} \right) & \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} - 2 \left(\frac{1}{\sigma^{2}} + \frac{2}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} - 2 \left(\frac{1}{\sigma^{2}} + \frac{2}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} - 2 \left(\frac{1}{\sigma^{2}} + \frac{2}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} - 2 \left(\frac{1}{\sigma^{2}} + \frac{2}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} - 2 \left(\frac{1}{\sigma^{2}} + \frac{2}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} - 2 \left(\frac{1}{\sigma^{2}} + \frac{2}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} - 2 \left(\frac{1}{\sigma^{2}} + \frac{2}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{\sigma^{2}} \right) \mu^{2} \\ = \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{2} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{2} \right) \mu^{2} + \frac{1}{2} \left(\frac{n}{\sigma^{2}} + \frac{1}{2} \right) \mu^{2} \\ = \frac{1}{2}$$

	3
For the multivariate case:	
p(x/u) ~ N(µ, E) 2 p(µ)~ N(µ0, 20)	
evoir Bayes formula, $P(\mu D) = \alpha \stackrel{?}{R} P(\alpha_{k} \mu) P(\mu)$	
116 to (5) in Ber -1[ut(nz1+z1)u-2ut) P(UID) = re	(5 5xx+ 5, 40)
	dan sa Wasii
	•

Solution 3: PCA, LDA and Gaussian Naive Bayes Classifier

PCA:

Principal component analysis is often used to reduce the dimensionality and bring out strong patterns in a dataset. The steps in PCA are:

- Centre the dataset, Compute the d-dimensional mean vector (i.e., the means for every dimension of the whole dataset)
- Compute the the covariance matrix of the whole data set
- Compute eigenvectors (e1,e2,...,ed) and corresponding eigenvalues (λ1,λ2,...,λd)
- Sort the eigenvectors by decreasing eigenvalues and choose k
 eigenvectors with the largest eigenvalues to form a d×k dimensional
 matrix W
- Use this d×k eigenvector matrix to transform the samples onto the new subspace. This can be summarized by the mathematical equation: y=W^T×x (where x is a d×1-dimensional vector representing one sample, and y is the transformed k×1-dimensional sample in the new subspace.)
- Since in the particular dataset, d >> n, Kernel PCA was used.

Now if A is an $m \times n$ matrix, with n >> m, then A^TA is a very large $n \times n$ matrix. So instead of computing the eigenvectors of A^TA , we might like to compute the eigenvectors of the much smaller $m \times m$ matrix AA^T -- assuming we can figure out a relationship between the two. So how are the eigenvectors of A^TA related to the eigenvectors of AA^T ?

Let v be an eigenvector of AA^T with eigenvalue λ . Then

- $AA^Tv = \lambda v$
- $A^T(AA^Tv) = A^T(\lambda v)$
- $(A^TA)(A^Tv) = \lambda(A^Tv)$

```
lines = open('dataset/dorothea_train.data', 'r').readlines()
dataList = [line.rstrip('\n') for line in lines]
k = len(max(dataList, key=len))
cols = np.arange(k)
train data = pd.read csv(
     filepath or buffer='dataset/dorothea train.data',
    header=None,
    sep=" ",
    names=cols,
    engine = 'python')
train data=train data.fillna(0)
X = lil matrix((train data.shape[0],100001))
for i in range(train data.shape[0]):
   X[i,train data[i]] = 1
k=500
K=X.dot(X.T)
kern cent = KernelCenterer()
S = kern cent.fit transform(K.toarray())
eig vals, eig vecs = np.linalg.eig(S)
eig pairs = [(np.abs(eig vals[i]), eig vecs[:,i]) for i in range(len(eig vals))]
# Sort the (eigenvalue, eigenvector) tuples from high to low
eig_pairs = sorted(eig_pairs, key=lambda k: k[0], reverse=True)
vec = np.array([ eig_pairs[i][1] for i in range(k)])
vec = vec.T # to make eigen vector matrix nxk
W=X.T.dot(vec)
Y=X.dot(W)
global pca_X
pca X = deepcopy(Y)
```

LDA:

- Computing the d-dimensional mean vectors
- Computing the Scatter Matrices
 - Within-class scatter matrix SW
 - Between-class scatter matrix SB
- Solving the generalized eigenvalue problem for the matrix S⁻¹W*SB
- Selecting linear discriminants for the new feature subspace
 - Choosing k eigenvectors with the largest eigenvalues
- Transforming the samples onto the new subspace

```
n = pca X.shape[0]
d = pca X.shape[1]
overall_mean = np.mean(pca_X, axis=0)
mean_vectors = []
for class_no in range(2):
    mean vectors.append(np.mean(pca X[class indices[class no]], axis=0))
SW = np.zeros((d,d))
for class_no,mean_vector in zip(range(2), mean_vectors):
    scatter_matrix = np.zeros((d,d))
    for row in pca X[class indices[class no]]:
        row, mean_vector = row.reshape(d,1), mean_vector.reshape(d,1)
scatter_matrix += (row-mean_vector).dot((row-mean_vector).T)
    S W += scatter matrix
SB = np.zeros((d,d))
for i,mean vector in enumerate(mean vectors):
    n = pca X[class indices[i],:].shape[0]
    mean vector = mean vector.reshape(d,1)
    overall mean = overall mean.reshape(d,1) # make column vector
    S B += n * (mean vector - overall mean).dot((mean vector - overall mean).T)
eig_vals, eig_vecs = np.linalg.eig(np.linalg.inv(S_W).dot(S_B))
eig_pairs = [(np.abs(eig_vals[i]), eig_vecs[:,i]) for i in range(len(eig_vals))]
eig pairs = sorted(eig_pairs, key=lambda k: k[0], reverse=True)
W = eig pairs[0][1]
global lda X
lda X = np.real(pca X.dot(W))
```

Gaussian Naive Bayes Classifier:

The Classifier model is defined as:

```
class Model:
    def __init__(self, X, labels):
        self.dataset = X
        self.class_labels = labels
        self.test_class_labels = list(labels)
        self.trainingSet= []
        self.testSet = []
        self.summaries = {}
```

Here summaries is a dictionary object that holds mean and standard deviations for every class.

First the data is divided as test and train and also seperated on the basis of class labels :

```
def k fold split(self,splitRatio):
    trainSize = int(len(self.dataset) * splitRatio)
    trainSet = []
    copy = list(self.dataset)
    while len(trainSet) < trainSize:</pre>
        index = random.randrange(len(copy))
        trainSet.append(copy.pop(index))
        self.test class labels.pop(index)
    self.trainingSet = trainSet
    self.testSet = copy
def seperate_by_class(self):
    separated = {}
    for i in range(len(self.dataset)):
        vector = self.dataset[i]
        label = self.class labels[i]
        if (label not in separated):
            separated[label] = []
        separated[label].append(vector)
    return separated
```

The classifier is then trained, calculation of summaries for each class.

```
def train_classifier(self,data_flag):
    separated = self.seperate_by_class()
    summaries = {}
    for classValue, instances in separated.items():
        if data_flag is 0:
            summaries[classValue] = self.summarize_pca(instances)
        elif data_flag is 1:
            summaries[classValue] = self.summarize_lda(instances)
        self.summaries = summaries
```

Testing phase: For predicting values for the test set, probabilities are calculated using:-

```
def test_classifier(self):
    predictions = []
    for i in range(len(self.testSet)):
        result = self.predict(self.testSet[i])
        predictions.append(result)
    return predictions

def get_accuracy(self,predictions):
    correct = 0
    for i in range(len(self.testSet)):
        if self.test_class_labels[i] == predictions[i]:
            correct += 1
    return (correct/float(len(self.testSet))) * 100.0
```

```
Results obtained are:
      Split of train and test data = 70 % (train)
       Test Run 1:
       K = 100
Size of train= 560 and test= 240
Accuracy after PCA(K==100):
                                       49.58333333333333
Size of train= 560 and test= 240
Accuracy after LDA: 79.58333333333333
       Test Run 2:
       K = 100
In [1]: runfile('/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3/prob2.py',
wdir='/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3')
Size of train= 560 and test= 240
Accuracy after PCA(K==100): 55.4166666666667
Size of train= 560 and test= 240
Accuracy after LDA: 77.91666666666667
      Test Run 1:
      K = 500
Size of train= 560 and test= 240
Accuracy after PCA(K==500): 88.33333333333333
Size of train= 560 and test= 240
Accuracy after LDA: 95.41666666666667
      Test Run 2:
      K = 500
In [7]: runfile('/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3/prob3.py',
wdir='/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3')
Reloaded modules: pca, lda
Exception ignored in: <_io.FileIO name='dataset/dorothea_train.data' mode='rb' closefd=True>
ResourceWarning: unclosed file < io.TextIOWrapper name='dataset/dorothea train.data' mode='r'
encoding='UTF-8'>
Size of train= 560 and test= 240
Accuracy after PCA(K==500): 88.75
Size of train= 560 and test= 240
Accuracy after LDA: 96.25
```

Observations:

The accuracy for PCA is low, since it doesn't take into consideration the classification of the data. However using LDA significantly increases accuracy since it considers maximises inter class distance and minimizes intra class distance.

Also as K increases, accuracy is increasing, since the number of dimensions retained for LDA is higher and hence more information is captured for classification.

On changing the split ratio for train vs test data, the accuracy does not significantly vary from previous results

In [8]: runfile('/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3/prob3.py',

```
wdir='/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3')
Reloaded modules: pca, lda
Exception ignored in: < io.FileIO name='dataset/dorothea train.data' mode='rb' closefd=True>
ResourceWarning: unclosed file < io.TextIOWrapper name='dataset/dorothea train.data' mode='r'
encoding='UTF-8'>
Size of train= 400 and test= 400
Accuracy after PCA(K==500): 90.25
Size of train= 400 and test= 400
Accuracy after LDA: 96.75
In [9]: runfile('/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3/prob3.py',
wdir='/home/levicorpus/Documents/UG3/SMAI/A2/Assign2/Prob3')
Reloaded modules: pca, lda
Exception ignored in: < io.FileIO name='dataset/dorothea train.data' mode='rb' closefd=True>
ResourceWarning: unclosed file < io.TextIOWrapper name='dataset/dorothea train.data' mode='r'
encoding='UTF-8'>
Size of train= 400 and test= 400
Accuracy after PCA(K==500): 90.25
Size of train= 400 and test= 400
Accuracy after LDA: 96.5
```