

SMAI Assignment 3 Report

Problem 1

Code

```
import numpy as np
import csv
import pdb
from collections import defaultdict

def nb_train(train_data,plabel,nlabel):
    num_attributes = len(train_data[0]) - 1
    pdict = [defaultdict(float) for x in range(num_attributes)]
    ndict = [defaultdict(float) for x in range(num_attributes)]
    plist = [0]*num_attributes
    nlist = [0]*num_attributes
    num_psample = 0
    num_nsample = 0
    for row in train_data:
        # count number of positive and negative labels
        class_label = row[-1]
        if class_label == plabel:
            num_psample += 1
            cdiclist = pdict
            clist = plist
        elif class_label == nlabel:
            num_nsample += 1
            cdiclist = ndict
            clist = nlist

        # compute feature freq and attr freq for each class
        feat_vec = row[:-1]
        ind = -1
        for feat in feat_vec:
            ind += 1
            if feat == '?':
                continue
            cdiclist[ind][feat] += 1
            clist[ind] += 1
```

```

    # compute priors
    priors = [np.log(float(num_psample)/(num_psample + num_nsample)),
np.log(float(num_nsample)/(num_psample + num_nsample))]
    # take log of probabilities
    for i in range(num_attributes):
        for feat in pdict[i]:
            if pdict[i][feat]!=0:
                pdict[i][feat] = np.log(pdict[i][feat]/plist[i])
            else:
                pdict[i][feat] = np.log(0.000000001)
        for feat in ndict[i]:
            if ndict[i][feat]!=0:
                ndict[i][feat] = np.log(ndict[i][feat]/nlist[i])
            else:
                ndict[i][feat] = np.log(0.000000001)
    return [priors, pdict, ndict]

```

```

def nb_predict(priors, pdict, ndict, test_data, plabel, nlabel):
    classified = 0
    misclassified = 0
    ans = 0

    for row in test_data:
        # compute posterior probability for each test sample against
each class
        ground_truth = row[-1]
        feat_vec = row[:-1]
        ind = -1
        pcumulative = priors[0]
        ncumulative = priors[1]
        for feat in feat_vec:
            ind += 1
            if feat=='?':
                continue
            pcumulative += pdict[ind][feat]
            ncumulative += ndict[ind][feat]
        # predict label based on posterior probabilities
        if pcumulative > ncumulative:
            predicted = plabel

```

```

    else:
        predicted = nlabel
        # check against ground truth
        if predicted == ground_truth:
            classified += 1
        else:
            misclassified += 1
        print "Accuracy",float(classified)/(classified+misclassified),"
over ", classified + misclassified, " samples."
    return classified, misclassified

total_classified = 0
total_misclassified = 0
# process data
raw_data = csv.reader(open("census-income.data"))
data0 = []
for i in list(raw_data):
    data0.append([item.strip() for item in i])
data = np.array(data0)
data = data[:,[2,3,4,41]]
np.random.shuffle(data)
data = data[:10000]
acc = []
for i in range(10):
    np.random.shuffle(data)
    train_data = data[0:5000]
    test_data = data[5001:10000]
    [priors, pdict, ndict] = nb_train(train_data,'50000+.',''-
50000.')
    [classified, misclassified] = nb_predict(priors, pdict, ndict,
test_data, '50000+.',''- 50000.')
    total_classified += classified
    total_misclassified += misclassified
    acc.append(float(classified)/(classified+misclassified))
acc = np.array(acc)
mean_acc =
float(total_classified)/(total_classified+total_misclassified)
sd = np.std(acc)
print "Mean Accuracy",mean_acc,"Standard Deviation",sd

```

Observations

```
joycode@nelovo:~/sem5/smai/Assign2/census$ python p1.py
Accuracy 0.813162632527 over 4999 samples.
Accuracy 0.823364672935 over 4999 samples.
Accuracy 0.828965793159 over 4999 samples.
Accuracy 0.831566313263 over 4999 samples.
Accuracy 0.831766353271 over 4999 samples.
Accuracy 0.806561312262 over 4999 samples.
Accuracy 0.831966393279 over 4999 samples.
Accuracy 0.806361272254 over 4999 samples.
Accuracy 0.801960392078 over 4999 samples.
Accuracy 0.835767153431 over 4999 samples.
Mean Accuracy 0.821144228846 Standard Deviation 0.0121713947933
```

- Ties are resolved by giving default class to class B
- Two things can be done for missing features:
 - Omit the records with missing values.
 - Omit only the missing attributes, where attributes are the values the feature can take. They are not considered while multiplying the probabilities, but are taken into account while general division.

Problem 2

Bayesian Parameter Estimation for Univariate Density Functions.

We know that class posterior probability can be estimated as:

$$p(w_i | x, D_i) = \frac{p(x | w_i, \theta_i) p(w_i)}{\sum_{j=1}^K p(x | w_j, \theta_j) p(w_j)}$$

where w_i : class

D_i : data samples for that class

fixing w_i : only 1 class, remove w_i

$p(x | \theta) = \int p(x | \theta) p(\theta | D) d\theta$. Let $\hat{\theta}$ is our optimal parameter.. θ (say mean).

$$p(x | D) \approx p(x | \hat{\theta}) \quad (\text{where } p(\theta | D) \text{ peaks}).$$

$$p(\theta | D) = \frac{p(\theta | D) p(\theta)}{\int_{\theta} p(\theta | D) p(\theta) d\theta}$$

We can think of this as following: Let the distribution $p(x | w_i, \theta)$ is a gaussian, with parameter, μ & σ . Then the parameter μ also varies like a gaussian

$$p(\mu | D) = \frac{p(D | \mu) p(\mu)}{\int_{\mu} p(D | \mu) p(\mu) d\mu} \quad - (1)$$

$p(\mu)$: prior: also a gaussian.

$$p(x | \mu) \sim \mathcal{N}(\mu, \sigma^2)$$

$$p(\mu) \sim \mathcal{N}(\mu_0, \sigma_0^2).$$

Consider (1)

$$p(\mu|D) = \frac{p(D|\mu) p(\mu)}{\int p(D|\mu) p(\mu) d\mu}$$

$$= \frac{1}{\int \prod_{k=1}^n p(x_k|\mu) p(\mu) d\mu} \quad \text{as all data points are independent.}$$

↳ constant normalizing term of Denominator

$$\text{As } p(x_k|\mu) \sim \mathcal{N}(\mu, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_k - \mu)^2}{2\sigma^2}}$$

$$\prod_{k=1}^n p(x_k|\mu) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{k=1}^n \frac{(x_k - \mu)^2}{2\sigma^2}}$$

$$p(\mu) = \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}}$$

$$\therefore p(\mu|D) \propto \exp\left(-\frac{1}{2} \left[\sum_{k=1}^n \frac{(x_k - \mu)^2}{\sigma^2} + \frac{(\mu - \mu_0)^2}{\sigma_0^2} \right]\right)$$

$$= \exp\left(-\frac{1}{2} \left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right) \mu^2 - 2\left(\frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right) \mu \right]\right)$$

$$\sigma_n^2 = \frac{\sigma^2 \sigma_0^2}{n\sigma_0^2 + \sigma^2}$$

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2}\right) \text{mean samples} + \left(\frac{\sigma^2}{n\sigma_0^2 + \sigma^2}\right) \mu_0$$

Bayesian Parameter Estimation for Multivariate Density Functions

Modifying the above result for a multivariate density function. \bar{x} : d-dimensional input, m such inputs

$$p(w_i | \bar{x}, D_i) = \frac{p(\bar{x} | w_i, b_i) p(w_i)}{\sum_{j=1}^J p(\bar{x} | w_j, b_j) p(w_j)}$$

$$p(\bar{x} | D) = \int p(\bar{x} | \bar{\theta}) p(\bar{\theta} | D)$$

where $\bar{\theta}$ is a vector of (multiple) parameters

approximating,

$$p(\bar{x} | D) \approx p(\bar{x} | \bar{\theta})$$

$$p(\bar{\theta} | D) = \frac{p(D | \bar{\theta}) p(\bar{\theta})}{\int_{\bar{\theta}} p(D | \bar{\theta}) p(\bar{\theta}) d\bar{\theta}}$$

$$p(D | \bar{\theta}) = \prod_{i=1}^m p(x_i | \bar{\theta})$$

$$p(\bar{\theta} | D) \propto \prod_{i=1}^m p(x_i | \bar{\theta}) p(\bar{\theta})$$

Also, $d \propto \frac{1}{\sqrt{\det(\Sigma)}}$

$$p(\bar{x}_i | \bar{\theta}) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left(-\frac{1}{2} (\bar{x}_i - \bar{\theta})^T \Sigma^{-1} (\bar{x}_i - \bar{\theta})\right) \sim \mathcal{N}(\bar{\theta}, \Sigma)$$

$$p(\bar{\theta}) = \frac{1}{(2\pi)^{d/2} \det(\Sigma_0)^{1/2}} \exp\left(-\frac{1}{2} (\bar{\theta} - \bar{\theta}_0)^T \Sigma_0^{-1} (\bar{\theta} - \bar{\theta}_0)\right)$$

$$\sim \mathcal{N}(\bar{\theta}_0, \Sigma_0)$$

$$\therefore p(\bar{\theta}|n) = d^n \exp(-\frac{1}{2})$$

$$\prod_{k=1}^n p(x_k|\bar{\theta}) = \frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{\det(Z_0)^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \sum_{k=1}^n \frac{(x_k - \bar{\theta})^T Z_0^{-1} (x_k - \bar{\theta})}{2}\right)$$

$$\therefore p(\bar{\theta}|0) = d^n \exp\left(-\frac{1}{2} \sum_{k=1}^n \frac{(x_k - \bar{\theta})^T Z_0^{-1} (x_k - \bar{\theta})}{2} - \frac{1}{2} (\bar{\theta} - \bar{\theta}_0)^T Z_0^{-1} (\bar{\theta} - \bar{\theta}_0)\right)$$

$$= d^n \exp\left(-\frac{1}{2} \left(\underbrace{\sum_{k=1}^n (x_k - \bar{\theta})^T Z_0^{-1} (x_k - \bar{\theta})}_A + \underbrace{(\bar{\theta} - \bar{\theta}_0)^T Z_0^{-1} (\bar{\theta} - \bar{\theta}_0)}_B \right)\right)$$

Consider A

$$= \sum_{k=1}^n (\bar{x}_k - \bar{\theta})^T Z_0^{-1} (\bar{x}_k - \bar{\theta}) = \sum_{k=1}^n (\bar{x}_k - \bar{\theta}_0 + \bar{\theta}_0 - \bar{\theta})^T Z_0^{-1} (\bar{x}_k - \bar{\theta}_0 + \bar{\theta}_0 - \bar{\theta}) = \underbrace{\sum_{k=1}^n (\bar{x}_k - \bar{\theta}_0)^T Z_0^{-1} (\bar{x}_k - \bar{\theta}_0)}_A + \underbrace{\sum_{k=1}^n (\bar{\theta}_0 - \bar{\theta})^T Z_0^{-1} (\bar{\theta}_0 - \bar{\theta})}_B + \underbrace{\sum_{k=1}^n (\bar{x}_k - \bar{\theta}_0)^T Z_0^{-1} (\bar{\theta}_0 - \bar{\theta})}_C$$

Consider B

$$= \underbrace{\bar{\theta}^T Z_0^{-1} \bar{\theta}}_A - \underbrace{\bar{\theta}^T Z_0^{-1} \bar{\theta}_0}_B - \underbrace{\bar{\theta}_0^T Z_0^{-1} \bar{\theta}}_C + \underbrace{\bar{\theta}_0^T Z_0^{-1} \bar{\theta}_0}_D$$

Combining A and B, we and completing the square in quadratic form, we obtain the final result as

$$\hat{\mu}_n = Z_0 (Z_0 + \frac{1}{n} S)^{-1} \hat{\mu}_n + \frac{1}{n} Z_0 (Z_0 + \frac{1}{n} S)^{-1} \bar{\theta}_0$$

$$\hat{\Sigma}_n = Z_0 (Z_0 + \frac{1}{n} S)^{-1} Z_0$$

where: $\hat{\mu}_n$ = mean of all samples

n = # samples

Problem 3

Code

```
clear;
clc;

fid = fopen('DOROTHEA/dorothea_train.data');
np = 800;
nvar = 100000;
X = zeros(nvar,np);
for i=1:np
    tline = fgetl(fid);
    ind = int32(str2double(strsplit(tline)));
    ind = ind(1:end-1);
    X(ind,i) = 1;
end
fclose(fid);
Y = load('DOROTHEA/dorothea_train.labels');

%PCA
k = 100;
X = X - mean(X,2)*ones(1,np);
K = X' * X;
[U,L] = eig(K);
[~,ind] = sort(diag(L),'descend');
PC = X * U(:,1:k);
X = PC' * X;

%LDA
% p = randperm(nvar, 10000);
% X = X(p,:);
% W1 = find(Y == -1);
% W2 = find(Y == 1);
% X_1 = X(:,W1);
% X_2 = X(:,W2);
% Xc_1 = X_1 - mean(X_1,2) * ones(1,size(W1,1));
% Xc_2 = X_2 - mean(X_2,2) * ones(1,size(W2,1));
% disp('done');
```

```

% Sw = Xc_1 * Xc_1' + Xc_2 * Xc_2';
% X = Sw \ (mean(X_1,2) - mean(X_2,2));

C = cvpartition(Y,'KFold',8);

meanacc = 0;

for epo=1:8

    tradat = X(:,training(C,epo));
    tralab = Y(training(C,epo));

    tesdat = X(:,test(C,epo));
    teslab = Y(test(C,epo));

    uniqlab = unique(tralab);
    nclass = length(uniqlab);
    nvar = size(X,1);
    ntest = length(teslab);

    for i=1:nclass
        ftralab(i) = sum(double(tralab==uniqlab(i)))/length(tralab);
    end

    for i=1:nclass
        tradat_i = tradat(:,(tralab==uniqlab(i)));
        tradat_mean(:,i) = mean(tradat_i,2);
        tradat_std(:,i) = std(tradat_i,0,2);
    end

    for i=1:ntest
        fteslab =
normcdf(ones(nclass,1)*tesdat(:,i)',tradat_mean',tradat_std');
        P(i,:) = ftralab.*prod(fteslab,2)';
    end

    [pv0,id]=max(P,[],2);
    for i=1:length(id)
        pv(i,1)=uniqlab(id(i));
    end

```

end

```
acc = sum(pv == teslab)/length(pv);  
meanacc = meanacc + acc;
```

end

```
disp(['accuracy = ', num2str(meanacc*(100/8)), '%'])
```

Observations

- Used k-fold cross validation on sample data; k = 8 i.e. test = 100, train = 700
- PCA k = 100, accuracy = 60.625%
- PCA k = 500, accuracy = 90.25%
- PCA k = 800, accuracy = accuracy = 90.25%