

```
In [1]: import numpy as np
import pandas as pd
```

```
In [2]: ## Import train data and the test data
X_train = np.load('hw2p2_train_x.npy')
X_test = np.load('hw2p2_test_x.npy')

y_train = np.load('hw2p2_train_y.npy')
y_test = np.load('hw2p2_test_y.npy')
```

Part(c)

(i)

```
In [21]: ## Get the train data where Y label is 1
index_y_is_1 = np.where(y_train == 1)[0]
X_label1 = X_train[index_y_is_1]
```

```
In [24]: ## Get the train data where Y label is 0
index_y_is_0 = np.where(y_train == 0)
X_label0 = X_train[index_y_is_0]
```

```
In [43]: alpha = 1
d = 1000
```

```
In [48]: ## Get a list of log(p_1j) for j = 1, ... 1000
n_k1 = np.sum(X_label1)
p_1j = []
for j in range(1000):
    frequency = 0
    for i in range(X_label1.shape[0]):
        frequency = frequency + X_label1[i][j]
    probs = (frequency + alpha)/(n_k1 + alpha * d)
    p_1j.append(np.log(probs))
```

```
In [52]: p_1j[:5]
```

```
Out[52]: [-7.024471078678098,
-7.717618259238043,
-7.247614629992308,
-7.717618259238043,
-7.38114602261683]
```

```
In [49]: ## Get a list of log(p_0j) for j = 1, ... 1000
n_k0 = np.sum(X_label0)
p_0j = []
for j in range(1000):
    frequency = 0
    for i in range(X_label0.shape[0]):
        frequency = frequency + X_label0[i][j]
    probs = (frequency + alpha)/(n_k0 + alpha * d)
    p_0j.append(np.log(probs))
```

```
In [53]: p_0j[:5]
```

```
Out[53]: [-6.055718936974995,
          -9.552226498441476,
          -9.552226498441476,
          -9.552226498441476,
          -9.552226498441476]
```

(ii)

```
In [71]: ## Compute the prior pi_0 and pi_1
estimate_pi_1 = np.log(X_label1.shape[0] / X_train.shape[0])
estimate_pi_0 = np.log(X_label0.shape[0] / X_train.shape[0])
```

```
In [72]: print("Estimate of prior pi_0 is:", estimate_pi_0, "Estimate of prior pi_1 is:",
```

```
Estimate of prior pi_0 is: -0.6965085282626502 Estimate of prior pi_1 is: -0.6897
970936746632
```

Part(d)

```
In [73]: prediction = []
for i in range(X_test.shape[0]):
    y0_value = 0
    y1_value = 0
    ### Get the value of belong to label 1 or label 0
    for j in range(1000):
        y0_value += X_test[i][j] * p_0j[j]
        y1_value += X_test[i][j] * p_1j[j]

    y0_value += estimate_pi_0
    y1_value += estimate_pi_1
    ### decision rule
    if y0_value > y1_value:
        prediction.append(0)
    else:
        prediction.append(1)
```

```
In [74]: ##define a function that to get the accuracy
def accuracy_bayes(X, Y):
    final_result = []
    for i in range(len(X)):
        if X[i] == Y[i]:
            final_result.append(1)
        else:
            final_result.append(0)
    return sum(final_result) / len(final_result)
```

```
In [75]: test_error = 1 - accuracy_bayes(prediction, y_test)
```

```
In [76]: print("The test error is for the naive bayesian classifier is:", test_error)
```

```
The test error is for the naive bayesian classifier is: 0.12594458438287148
```

Part(e)

```
In [77]: if_list = [1] * X_test.shape[0]
```

```
In [78]: if_test_error = 1 - accuracy_bayes(if_list, y_test)
```

```
In [79]: print("The test error is for the naive bayesian classifier is:", if_test_error)
```

The test error is for the naive bayesian classifier is: 0.49874055415617125