# ECE4950 HW2

Yi Chen yc2329 https://github.com/agraynel

March 7, 2017

# 1 Problem1

## 1.1 class probabilities

$P_r(Gender = M) = \frac{4}{9}$

$P_r(Gender = F) = \frac{5}{9}$

## 1.2 Estimate means and variances

Means:

$\mu_{\mathcal{H}|\mathcal{M}} = (72 + 68 + 75 + 64) \ / \ 4 = 69.75$

$\mu_{\mathcal{S}|\mathcal{M}} = (12 + 9 + 11 + 10.5) \ / \ 4 = 10.625$

$\mu_{\mathcal{H}|\mathcal{F}} = (65 + 67 + 62 + 70 + 64) \ / \ 5 = 65.6$

$\mu_{\mathcal{S}|\mathcal{F}} = (8 + 7.5 + 6 + 8.5 + 8) \ / \ 5 = 7.6$

Variances:

$\sigma^2_{\mathcal{H}|\mathcal{M}} = \frac{1}{N} \sum_{i=1}^{N} (x(i) - \mu)^2$

$= ((72 - 69.75)^2 + (68 - 69.75)^2 + (75 - 69.75)^2 + (64 - 69.75)^2)/4$

$= 17.1875$

$$\sigma^2_{\mathcal{S}|\mathcal{M}} = ((12 - 10.625)^2 + (9 - 10.625)^2 + (11 - 10.625)^2 + (10.5 - 10.625)^2)/4$$
$$= 1.172$$
$$\sigma^2_{\mathcal{H}|\mathcal{F}} = ((65 - 65.6)^2 + (67 - 65.6)^2 + (62 - 65.6)^2 + (70 - 65.6)^2 + (64 - 65.6)^2)/5$$
$$= 7.44$$
$$\sigma^2_{\mathcal{S}|\mathcal{F}} = ((8 - 7.6)^2 + (7.5 - 7.6)^2 + (6 - 7.6)^2 + (8.5 - 7.6)^2 + (8 - 7.6)^2)/5 = 0.74$$

## 1.3   Gausian Naive Bayes

Suppose height is 68 inches, and foot size is 9.5:

$$p(H = 68|M) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(v - \mu_c)^2}{2\sigma_c^2}} = \frac{1}{10.39} e^{-\frac{3.0625}{34.375}} = 0.088$$

$$p(S = 9.5|M) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(v - \mu_c)^2}{2\sigma_c^2}} = \frac{1}{2.714} e^{-\frac{1.266}{2.344}} = 0.215$$

$$p(H = 68|F) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(v - \mu_c)^2}{2\sigma_c^2}} = \frac{1}{6.837} e^{-\frac{5.76}{14.88}} = 0.099$$

$$p(S = 9.5|F) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(v - \mu_c)^2}{2\sigma_c^2}} = \frac{1}{2.156} e^{-\frac{3.61}{1.48}} = 0.040$$

## 1.4   Classification

$P_r(Gender = M) = \frac{4}{9}$

$P_r(Gender = F) = \frac{5}{9}$

$p(H, S|M) = p(H|M) * p(S|M) = 0.01892$

$p(M)p(H, S|M) = \frac{4}{9} * 0.01892 = 0.00841$

$p(H, S|F) = p(H|F) * p(S|F) = 0.00396$

$p(F)p(H, S|F) = \frac{5}{9} * 0.00396 = 0.0022$

Because $p(M)p(H, S|class = M)$ is larger, so we would classify the above sample as male.

# 2   Problem2

## 2.1   Gaussian Naive Bayes

See appendix and attached files for codes.
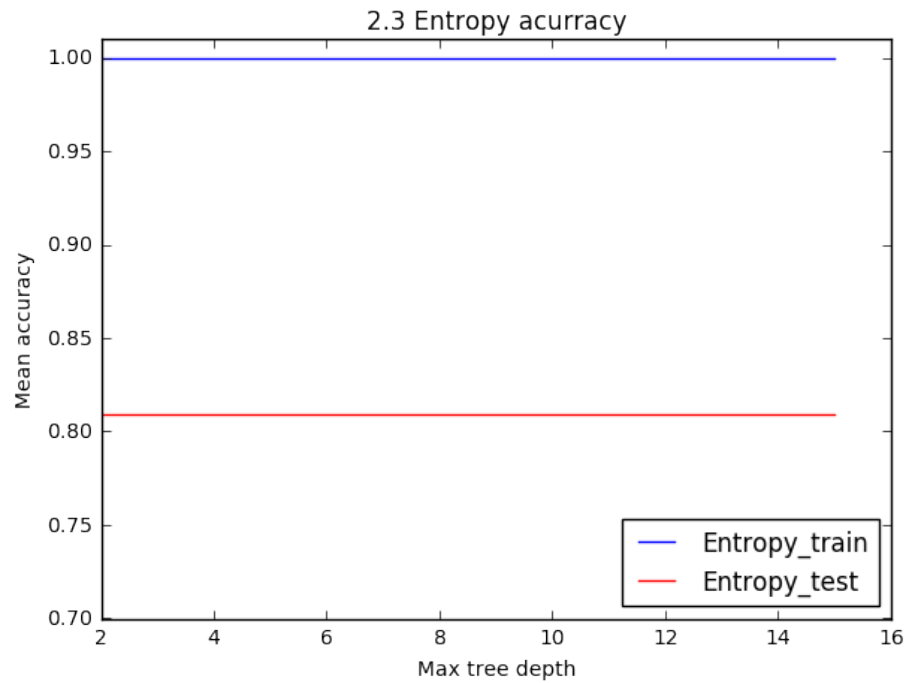
## 2.2  Accuracy

-50.csv: 0.890173410405

-200.csv: 0.964769647696

-400.csv: 0.96449704142

## 2.3  Decision tree learning

See appendix and attached files for codes.

The figure of accuracy of -50.csv using decision tree learning methods:



The test accuracy of decision tree with best depth is around 0.81 to 0.82, while the accuracy of Naive Bayes is 0.89. Obviously, Naive Bayes has a better test error.

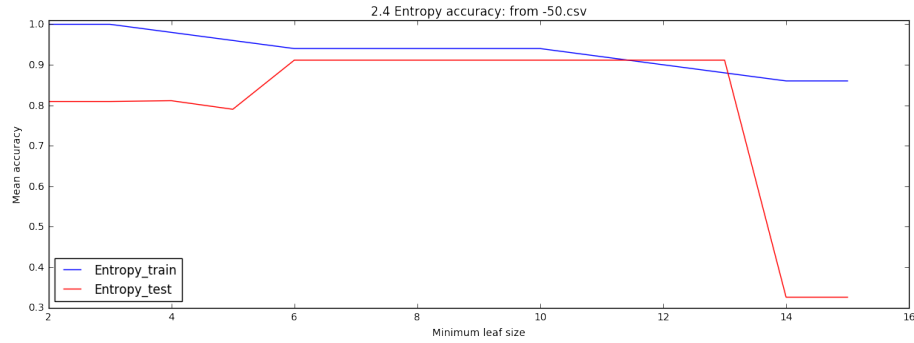## 2.4  Decision tree using minimum leaf size criterion

See appendix and attached files for codes, and figures are at the end of this section.

In -50.csv, compared to maximum depth criterion, when minimum leaf size is between 6 and 13(both sides included), the accuracy is larger, reaching 0.9. While the accuracy of decision tree classifier using maximum depth criterion is always above 0.8. When minimum leaf size is between 6 and 13, the accuracy is also a little larger than the accuracy of Naive Bayes.
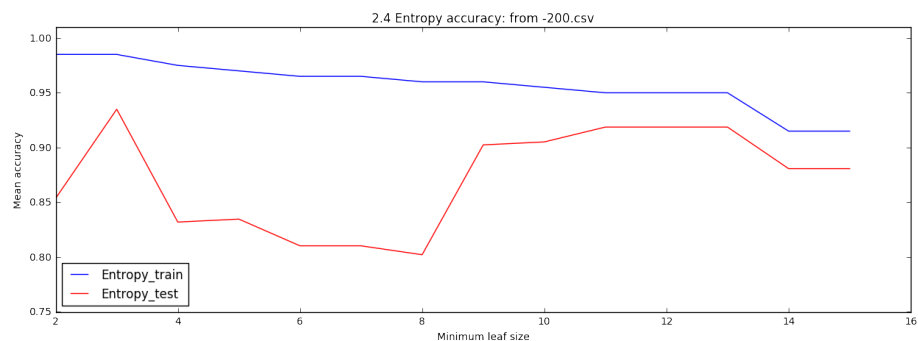
In -200.csv, compared to maximum depth criterion, when minimum leaf size is 2, 3 or between 9 and 13, the accuracy is larger. While the accuracy of decision tree classifier using maximum depth criterion is always around 0.8 when the depth is larger than 3. The accuracy is less than the accuracy of Naive Bayes.

In -400.csv, compared to maximum depth criterion, the maximum and average accuracy are smaller in using minimum leaf size. The accuracy is also less than the accuracy of Naive Bayes.
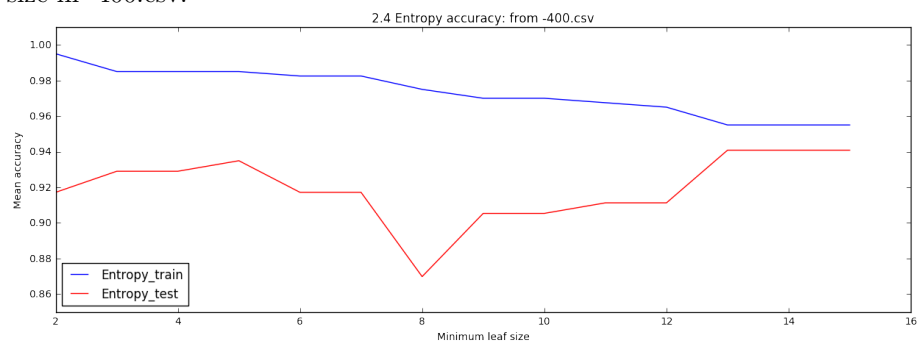
The figure of test and training accuracy with respect to different minimum leaf size in -50.csv:



The figure of test and training accuracy with respect to different minimum leaf size in -200.csv:

2.4 Entropy accuracy: from -200.csv

The figure of test and training accuracy with respect to different minimum leaf size in -400.csv:


2.4 Entropy accuracy: from -400.csv

From homework 1, we have the figure of test and training accuracy with respect to maximum depth in -200.csv:


4.2(b) Entropy acurracy


4.2(c) Gini acurracy

From homework 1, we have the figure of test and training accuracy with respect to maximum depth in -400.csv:

# 3 Problem3

## 3.1 Multinomial Naive Bayes

See appendix and attached files for codes.

Error rate and smoothness parameter:

| $\beta$ | Train error rate | Test error rate |
|---|---|---|
| 1e-07 | 0.174333333333 | 0.163 |
| 0.0001 | 0.17445 | 0.1631 |
| 0.1 | 0.174683333333 | 0.1633 |
| 1 | 0.174716666667 | 0.1635 |
| 1000 | 0.176466666667 | 0.1644 |
| 1000000 | 0.249066666667 | 0.2426 |
| 1000000000 | 0.26525 | 0.2616 |
| 1000000000000 | 0.887633333333 | 0.8865 |

Figure of error rate as a function of smoothness parameter $\beta$:

3.1 Error rate on smoothness parameter

## 3.2 $\beta \to \infty$

See appendix and attached files for codes.

$Add - \beta Prob(x) = \frac{n_x + \beta}{N + D*\beta}$

As $\beta \to \infty$, error rate also $\to 0.9. When smoothing, some of the features will not present in the the the learning samples$

| $\beta$ | Train error rate | Test error rate |
|---|---|---|
| inf | 0.901283333333 | 0.902 |

# 4  Appendix

## 4.1  Code for 2.1

```
1  #2.1 Gaussian Naive Bayes
2  #This is the code for all the 3 sets of data. Only need
     to change the file names in the code! On piazza, TA
```

```
       said  only  using  entropy  criterion  is  enough .
 3  #CREDITS:  http :// scikit −learn . org / stable /modules/
       generated / sklearn . naive_bayes . GaussianNB . html#sklearn .
       naive_bayes . GaussianNB . partial_fit

 4
 5  import  matplotlib . pyplot  as  plt
 6  from  sklearn . naive_bayes  import  GaussianNB
 7  import  pandas  as  pd
 8  import  numpy  as  np
 9  from  sklearn . metrics  import  accuracy_score

10
11  #this  is  the  path  of  dataset
12  Xtrn = pd . read_csv ( ' / Users /Agraynel /Desktop /ECE4950/
       homework/hw2/ dataset /X−trn −50. csv ' ,  header  =  None)
13  Xtst = pd . read_csv ( ' / Users /Agraynel /Desktop /ECE4950/
       homework/hw2/ dataset /X−tst −50. csv ' ,  header  =  None)
14  Ytrn = pd . read_csv ( ' / Users /Agraynel /Desktop /ECE4950/
       homework/hw2/ dataset /Y−trn −50. csv ' ,  header  =  None)
15  Ytst = pd . read_csv ( ' / Users /Agraynel /Desktop /ECE4950/
       homework/hw2/ dataset /Y−tst −50. csv ' ,  header  =  None)
16  #NOTICE:  change  the  file  name  to  −50.csv ,  −200.csv ,  −400.
       csv  respectively .

17
18  clf = GaussianNB ( )
19  clf . fit ( Xtrn ,  Ytrn . values . ravel ( ) )
20  #CREDITS:  http :// stackoverflow . com/ questions /34165731/a−
       column−vector −y−was−passed−when−a−1d−array −was−
       expected
21  #GaussianNB( priors=None)
22  #  train  with  Gaussian
```

```
23  test = clf.predict(Xtst)
24  train_accuracy = clf.score(Xtrn, Ytrn)
25  test_accuracy = accuracy_score(Ytst, test)
26
27  print(test_accuracy)
```

## 4.2    Code for 2.3

```
1  #2.3 read 50 files decision tree classifier with maximum
       depth
2
3  import matplotlib.pyplot as plt
4  from sklearn.tree import DecisionTreeClassifier
5  import pandas as pd
6  import numpy as np
7  from sklearn.metrics import accuracy_score
8
9  #this is the path of dataset
10 Xtrn50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/X-trn-50.csv', header = None)
11 Xtst50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/X-tst-50.csv', header = None)
12 Ytrn50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/Y-trn-50.csv', header = None)
13 Ytst50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/Y-tst-50.csv', header = None)
14
15 entropy_50_train = []
16 entropy_50_test = []
17
```

```python
18  for i in range(2, 16):
19      clf_entropy = DecisionTreeClassifier(criterion = '
            entropy', max_depth= i, random_state = 0)
20      clf_entropy = clf_entropy.fit(Xtrn50, Ytrn50)
21      entropy_test = clf_entropy.predict(Xtst50)
22      entropy_50_train.append(clf_entropy.score(Xtrn50,
            Ytrn50))
23      entropy_50_test.append(accuracy_score(Ytst50,
            entropy_test))
24
25  # Plotting decision regions
26
27  plt.figure(figsize=(15, 5))
28  plt.plot(range(2, 16), entropy_50_train, c='blue', label=
        'Entropy_train')
29  plt.plot(range(2, 16), entropy_50_test, c='red', label='
        Entropy_test')
30  plt.legend(loc=4)
31  plt.ylim(0.7, 1.01)
32  plt.ylabel('Mean_accuracy')
33  plt.xlabel('Max_tree_depth')
34  plt.title('2.3_Entropy_accuracy')
35
36  plt.show()
```

## 4.3   Code for 2.4

```python
1  #2.4 decision tree classifier with minimum leaf size i
       criterion.
```

```python
2  #this is one typical code for one situation. jusr change
       the path and get other datasets.
3
4  import matplotlib.pyplot as plt
5  from sklearn.tree import DecisionTreeClassifier
6  import pandas as pd
7  import numpy as np
8  from sklearn.metrics import accuracy_score
9
10 #this is the path of dataset
11 Xtrn50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/X-trn-50.csv', header = None)
12 Xtst50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/X-tst-50.csv', header = None)
13 Ytrn50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/Y-trn-50.csv', header = None)
14 Ytst50 = pd.read_csv('/Users/Agraynel/Desktop/ECE4950/
       homework/hw2/dataset/Y-tst-50.csv', header = None)
15
16 entropy_50_train = []
17 entropy_50_test = []
18 #change 50 to 200, 400, then we can acquire all the
       figures of 3 conditions.
19
20 for i in range(2, 16):
21     clf_entropy = DecisionTreeClassifier(criterion = '
           entropy', min_samples_leaf = i, random_state = 0)
22     clf_entropy = clf_entropy.fit(Xtrn50, Ytrn50)
23     entropy_test = clf_entropy.predict(Xtst50)
```

```
24        entropy_50_train.append(clf_entropy.score(Xtrn50,
              Ytrn50))
25        entropy_50_test.append(accuracy_score(Ytst50,
              entropy_test))
26
27 # Plotting decision regions
28
29 plt.figure(figsize=(15, 5))
30 plt.plot(range(2, 16), entropy_50_train, c='blue', label=
       'Entropy_train')
31 plt.plot(range(2, 16), entropy_50_test, c='red', label='
       Entropy_test')
32 plt.legend(loc = 3)
33 plt.ylim(0.3, 1.01)
34 plt.ylabel('Mean_accuracy')
35 plt.xlabel('Minimum_leaf_size')
36 plt.title('2.4_Entropy_accuracy')
37
38 plt.show()
```

## 4.4   Code for 3.1

```
1 #3.1 Error rate on smoothness parameter
2 import csv
3 import matplotlib.pyplot as plt
4 from sklearn.naive_bayes import MultinomialNB
5 import pandas as pd
6 import numpy as np
7 from sklearn.metrics import accuracy_score
8 from sklearn import datasets
```

```python
from sklearn.datasets import fetch_mldata

beta = [0.0000001, 0.0001, 0.1, 1, 1000, 1000000,
        1000000000, 1000000000000]
## Loading data
mnist = fetch_mldata('MNIST-original')
train, test = mnist.data[0:60000, :], mnist.data[60000:,
        :]
X_trn, y_trn = mnist.data[0:60000, :], mnist.target
        [0:60000]
X_tst, y_tst = mnist.data[60000:, :], mnist.target
        [60000:]

trainError = [];
testError = [];
for i in range(0, 8):
    clf = MultinomialNB(alpha = beta[i], class_prior =
            None, fit_prior = True)
    clf = clf.fit(X_trn, y_trn)
    tst = clf.predict(X_tst)
    accuracy = clf.score(X_trn, y_trn)
    testAccuracy = accuracy_score(y_tst, tst)
    trainError.append(1 - accuracy)
    testError.append(1 - testAccuracy)
    print(beta[i], 1 - accuracy, 1 - testAccuracy)

# Plotting error rate as a function of smoothness
        parameter

plt.figure()
```

```
33  plt.plot(beta, trainError, c='blue', label='trainError')
34  plt.plot(beta, testError, c='red', label='testError')
35  plt.xscale('log')
36  plt.legend(loc = 0)
37  plt.ylim(0.1, 1)
38  plt.ylabel('Error_rate')
39  plt.xlabel('Smoothness_parameter')
40  plt.title('3.1_Error_rate_on_smoothness_parameter')
41
42  plt.show()
```

## 4.5   Code for 3.2

```
1   #3.2 smoothness parameter is infinity
2   import csv
3   import matplotlib.pyplot as plt
4   from sklearn.naive_bayes import MultinomialNB
5   import pandas as pd
6   import numpy as np
7   from sklearn.metrics import accuracy_score
8   from sklearn import datasets
9   from sklearn.datasets import fetch_mldata
10
11  ## Loading data
12  mnist = fetch_mldata('MNIST-original')
13  train, test = mnist.data[0:60000, :], mnist.data[60000:,
        :]
14  X_trn, y_trn = mnist.data[0:60000, :], mnist.target
        [0:60000]
```

```python
15  X_tst, y_tst = mnist.data[60000:, :], mnist.target
        [60000:]

16

17  trainError = [];
18  testError = [];
19  clf = MultinomialNB(alpha = float("inf"), class_prior =
        None, fit_prior = True)
20  clf = clf.fit(X_trn, y_trn)
21  tst = clf.predict(X_tst)
22  accuracy = clf.score(X_trn, y_trn)
23  testAccuracy = accuracy_score(y_tst, tst)
24  trainError.append(1 - accuracy)
25  testError.append(1 - testAccuracy)
26  print(float("inf"), 1 - accuracy, 1 - testAccuracy)
```