

Total 4 sets of data have been downloaded below. Two sets for training, two sets for testing.

Each sample in set has 28 x 28 matrix to represent digit 0 or digit 1. Value of 0 in matrix means white, value of 255 is black. Color is grayscale and can be represented in between 0 and 255.

I need to extract two features for each image. Figure 1 is average all pixel brightness values within a whole image array. Feature 2 standard deviation of all pixel brightness values within a whole image array. Since I am going to use the Gaussian probability distribution function (GPDF) to classify digit 0 or digit 1, we need mean and standard deviation. I need to calculate 8 the parameters for the two-class (digit 0 and 1) naive bayes classifiers with GPDF.

Now, I will need calculate probability to differentiate digit 0 from digit 1. If probability of digit 0 is greater than probability of digit 1, I classify as digit 0 by using Naïve Bayse Classifier.

Here are formulas which I used for Naïve Bayse Classifier

$$P(y|x_1,x_2) = P(y)P(x_1|y) P(x_2|y) / P(x_1)P(x_2), \text{ ignore } P(x_1)P(x_2)$$
$$P(0|f_1,f_2) = P(0)P(f_1|0) P(f_2|0)$$
$$P(1|f_1,f_2) = P(1)P(f_1|1) P(f_2|1)$$

,where $P(0)=P(1)=0.5$

$$\begin{aligned} P(0|f_1, f_2) &> P(1|f_1, f_2), \text{ then digit} = 0 \\ P(0|f_1, f_2) &< P(1|f_1, f_2), \text{ then digit} = 1 \end{aligned}$$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Gaussian probability distribution function

```
# Calculate the Gaussian probability distribution function for x
def cal_gauss_prob(x, mean, var):
    return (1 / (sqrt(2 * pi * var))) * exp(-((x-mean)**2 / (2 * var)))
```

Python code for Gaussian probability distribution function

```
# predicting label for test 0 data
i=0 #predicting label digit 0 in test data set
t=0 #total number of test data set
for test in test0_features:
    #print(test[0], test[1])
    prob_0 = 0.5 * cal_gauss_prob(test[0], train0_means[0], train0_vars[0]) * cal_gauss_prob(test[1], train0_means[1], train0_vars[1])
    prob_1 = 0.5 * cal_gauss_prob(test[0], train1_means[0], train1_vars[0]) * cal_gauss_prob(test[1], train1_means[1], train1_vars[1])
    #print(prob_0, prob_1)
    if prob_0 > prob_1:
        i += 1
    t += 1
#print(i, t)
print('Accuracy_for_digit0testset = ', i/t)
```

Python code to calculate the accuracy of your predictions for testset Accuracy (digit 0 and 1).

After comparing probability of digit 0 with probability of digit 1, I picked up the label digit which has higher probability and count it in. Then, I am able to calculate accuracy by using Counts divided by total number as you could see at python code.

Here are all results of Project 1

```
Mean_of_feature1_for_digit0=      44.2838607143
Variance_of_feature1_for_digit0 =  114.410875544
Mean_of_feature2_for_digit0 =     87.5115913698
Variance_of_feature2_for_digit0 =  100.356180612
Mean_of_feature1_for_digit1 =     19.335694898
Variance_of_feature1_for_digit1 =   31.346498039
Mean_of_feature2_for_digit1 =     61.2965618747
Variance_of_feature2_for_digit1 =   82.046224968
Accuracy_for_digit0testset = 0.9173469387755102
Accuracy_for_digit1testset = 0.9233480176211454
```

```
[44.2838607143, 114.410875544, 87.5115913698, 100.356180612, 19.335694898, 31.346498039, 61.2965618747, 82.046224968, 0.9173469387755102, 0.9233480176211454]
```